

Exhibit 2

PLAINTIFFS' OPPOSITION TO DEFENDANTS' MOTION TO EXCLUDE PLAINTIFFS' EXPERTS' GENERAL CAUSATION OPINIONS FOR FAILURE TO ACCOUNT FOR SECTION 230 AND THE FIRST AMENDMENT

Case No.: 4:22-md-03047-YGR
MDL No. 3047

In Re: Social Media Adolescent Addiction/Personal Injury Products Liability Litigation

Highly Confidential (Competitor)

**SUPERIOR COURT OF THE STATE OF CALIFORNIA
FOR THE COUNTY OF LOS ANGELES**

COORDINATION PROCEEDING SPECIAL TITLE [RULE 3.400]

JUDICIAL COUNCIL COORDINATED
PROCEEDING NO. 5255

SOCIAL MEDIA CASES

Judge: Hon. Carolyn B. Kuhl

EXPERT REPORT OF DR. EMILIO FERRARA

**THIS DOCUMENT RELATES TO:
ALL CASES**

Highly Confidential (Competitor)

TABLE OF CONTENTS

I.	Introduction and Summary of Qualifications.....	5
A.	Educational Background	5
B.	Professional Experience	5
C.	Publications	12
II.	Plaintiffs' Allegations and Summary of Opinions	12
A.	Plaintiffs' Allegations Regarding Content Moderation	12
B.	Nature of Assignment and Summary of Opinions.....	12
III.	Background	14
A.	Overview of Content Moderation on Social Media Services	14
B.	Challenges to Content Moderation on Social Media Services	16
1.	Socio-political Challenges	17
a.	What harm is the most pressing is always evolving	17
b.	Ever-evolving regulatory challenges	19
2.	Technical Challenges for State of the Art Artificial Intelligence Systems	22
3.	Limitations of human reviewers.....	34
IV.	Meta's Content Moderation Policies Are Comprehensive and Appropriately Balanced.....	35
A.	Overview and Development of Content Moderation Policy	35
B.	Summary of Public-Facing Content Moderation Policies on Facebook or Instagram	40
1.	Suicide & Self-Injury & Eating Disorder Content	40
2.	Bullying and Harassment Content.....	43
3.	Adult Nudity & Sexual Activity Content	45
4.	Violent & Graphic Content.....	47
5.	Child Safety Content	49

Highly Confidential (Competitor)

1	C.	Borderline Policies.....	51
2	V.	Meta's Content Moderation Systems Are Robust and Reasonably Designed.....	52
3	A.	Proactive Detection.....	53
4	1.	Meta's Early Leadership in Proactive Detection Technology	53
5	a.	Facebook Immune System	53
6	b.	Research-Driven Technology.....	55
7	2.	How Classifiers Work	56
8	a.	Content Moderation of Live (Real-Time) Content.....	60
9	b.	Abusive Account Moderation.....	61
10	c.	Content-Specific Classifiers	62
11	d.	Suicide and Self-Injury Classifiers	62
12	e.	Eating Disorder Classifiers.....	64
13	3.	Bullying and Harassment Classifiers.....	65
14	4.	Graphic & Violent Content Classifiers.....	66
15	5.	Adult Nudity & Sexual Activity Classifiers	67
16	6.	Child Safety Classifiers	67
17	a.	Borderline Classifiers.....	69
18	b.	Content Downranking.....	69
19	c.	User Warning Screens and Labels	70
20	B.	Additional Automated Tools to Prevent Exposure to Harmful Content.....	70
21	1.	Public Efforts to Improve Content Moderation Technology	73
22	2.	Other Efforts to Improve Content Moderation Technology	75
23	C.	Meta's Reporting System	75
24	D.	Meta's Human Review Processes.....	77
25	1.	Overview.....	77
26	2.	In-House & Third-Party Human Reviewers	79
27	3.	Prioritization of Content for Human Review.....	80
28			

Highly Confidential (Competitor)

VI.	Effectiveness and Transparency of Meta's Content Moderation Systems	83
A.	CSERs & the Prevalence Metric.....	83
B.	Reasonableness of CSERs for Evaluating the Effectiveness of Content Moderation System	83
C.	Most-Recent CSER Report.....	85
D.	Expert Assessment of the Algorithmic Stress Test of Instagram's Reels Surface.....	88
1.	Mr. Bejar Conducted a Truncated Algorithmic Stress Test.....	89
2.	Mr. Bejar's Algorithmic Stress Test Suffers from Fatal Methodological Flaws Rendering the Results Scientifically Unreliable.....	92
VII.	Meta's Content Moderation Policies & Enforcement Align with Industry Best Practices	97
A.	Product Development: Proactive Risk Identification and Mitigation.....	98
B.	Product Governance: Transparent and Evolving Policy Structures.....	98
C.	Enforcement Operations: Scalable and Responsive Moderation Infrastructure	98
D.	Iterative Improvement: Data-Driven Learning and Policy Adjustment.....	98
E.	Transparency and Accountability: Industry-Leading Disclosures	99
VIII.	Conclusions	99
	Appendix A: Table of Harm-Specific Classifiers.....	101

Highly Confidential (Competitor)

1 **I. Introduction and Summary of Qualifications**

2 1. My name is Emilio Ferrara. I am a Professor of Computer Science at the University
3 of Southern California, Associate Chair at the Thomas Lord Department of Computer Science, and
4 the inaugural (ad interim) Director of the Interdisciplinary Data Science Program at the USC
5 School of Advanced Computing. I have joint appointments in Communication (USC Annenberg)
6 and Preventive Medicine (USC Keck School of Medicine). I am also Principal Scientist at the
7 Information Sciences Institute, Principal Investigator at the USC-ISI Machine Intelligence and
8 Data Science (MINDS) center, and Director of the PhD program in Computer Science.

9 2. I have summarized in this section my educational background, career history,
10 publications, and other relevant qualifications. A copy of my CV containing additional details
11 relating to my background and qualifications is attached as Exhibit A to this report.

12 **A. Educational Background**

13 3. I attended the University of Messina (Italy), where I received a bachelor's degree
14 in Computer Science in 2006, and a Master's degree in Computer Science in 2008. I completed
15 my doctoral studies by the end of 2011, which focused on Machine Learning, with a concentration
16 on the study of online social networks. My doctoral thesis titled "Mining and analysis of online
17 social networks" was successfully defended in March 2012. My PhD degree was awarded the
18 special honorific title of Doctor Europaeus, which indicates that the PhD student has met specific
19 requirements of the European University Association related to international collaboration,
20 mobility, and multilingualism.¹ During my graduate studies, I published 13 peer-reviewed research
21 papers. I have since published over 200 additional peer-reviewed research papers.

22 **B. Professional Experience**

23 4. In 2010, I served as a visiting scholar at Technische Universität Wien in Vienna,
24 Austria, where I was involved with research focusing on Data Mining and Machine Learning.

25 ¹ See *Doctor Europaeus*. (n.d.). Università degli Studi di Udine.

26 https://www.uniud.it/en/research/do-research/doctorate-res/international-programmes/doctor-
27 europaeus#:~:text=The%20Doctor%20Europaeus%20certification%20is,qualification%20issued
%20by%20international%20institutions.

Highly Confidential (Competitor)

1 While there, I worked as an intern at Lixto Software GmbH (later acquired by McKinsey), where
2 I was involved with research to develop Web crawler bots for large-scale, distributed data
3 collection from online services. This work appeared in 3 peer-reviewed research papers.

4 5. From 2011-2012, I was a visiting scholar at Royal Holloway, University of London
5 in Egham, England, where I was involved with research focusing on Network Science,
6 Computational Biology, and Machine Learning. During this period, I developed techniques for the
7 clustering of complex network data and applied these methodologies to social and biological
8 networks. This research has appeared in over a dozen peer-reviewed research articles. I
9 implemented these techniques and released them as open-source software, now used by hundreds
10 of research papers spanning applications in science, engineering, and business.

11 6. At Indiana University Bloomington (2012-2015), I conducted research focusing on
12 Network Science, Data Science, Data Mining, and Machine Learning. My roles included
13 postdoctoral research fellow, assistant research scientist, and research assistant professor within
14 the School of Engineering and Computing and the Indiana University Networks Institute (IUNI).
15 During this period, I published approximately 40 peer-reviewed research papers, which laid the
16 foundation for the study of Twitter's content moderation systems. My focus was the detection of
17 manipulation and abuse on the Twitter service. I contributed to the ideation and design, as a co-
18 inventor, of the "Bot Or Not" (later renamed "Botometer") framework for the detection of bots on
19 Twitter. I used this and other machine learning techniques I invented to study a multitude of social
20 phenomena on Twitter, including but not limited to (i) the coordination of Twitter users revolving
21 around the "Occupy Wall Street" social movements of 2012; (ii) the dynamics of evolution of
22 Twitter trends; (iii) the behavioral evolution of Twitter users revolving around the "Gezi Park"
23 protests of 2013 in Turkey; (iv) the issue of misogynistic language and toxic abuse of Twitter
24 users; (v) the issue of emotional contagion and the dynamic of spread of negative sentiment on
25 Twitter; (vi) the issue of manipulation and abuse, specifically the manufacturing of alarming
26 misinformation, on Twitter in the context of the 2014 Ebola public health crisis.

27 7. In 2015, I began working at the University of Southern California as a computer
28

Highly Confidential (Competitor)

1 scientist, conducting research into Network Science, Data Science, Data Mining, and Machine
2 Learning. In 2016, I transitioned to the role of research assistant professor at USC, and in 2020, I
3 became a tenured associate professor. In 2022, I was promoted to full professor, the same role I
4 have today. I currently hold appointments as professor of computer science, communication, and
5 public health, as well as associate chair of the Department of Computer Science.

6 8. In 2021, I also began working with Amazon as a visiting academic in their Alexa
7 AI department. In this role, I drive research initiatives to prevent Alexa from providing
8 inappropriate responses to customers. I do so by collaborating with Amazon's Alexa team in
9 charting the taxonomy of types of problematic content (e.g., hate speech, spam and scams, health
10 misinformation, adult content, violent content, racially-charged or other types of biased content,
11 and more) that should not be digested by, consumed through, or exposed to Alexa AI. I also design
12 machine learning techniques to (i) prevent indexing problematic content from the Web and social
13 media services; (ii) detect and filter out problematic content early on from the indexed data; (iii)
14 detect questions that are unsuitable to a machine-generated answer; (iv) detect unsuitable answers
15 to specific audiences (e.g., children); (v) detect social media trends that might be reflected in
16 inquiries to Alexa; (vi) anticipate inappropriate content from social media feeds.

17 9. Since 2015, I have developed and used a plethora of machine learning and AI
18 techniques to study a multitude of problems on social media services, including but not limited to
19 (i) detecting online extremism and radical propaganda by the ISIS terrorist group on social media;
20 (ii) detecting spam in electronic-cigarette advertising campaigns on social media; (iii) unveiling
21 manipulation of social media discourse revolving around the 2016 US Presidential Election; (iv)
22 automatically detecting bots by means of an ensemble of classifiers, known as "Bot Or Not" and
23 later "Botometer" framework; (v) participating in the Defense Advanced Research Projects
24 Agency social media bot challenge and ranking in the top 3 teams for accuracy of bot detection;
25 (vi) characterizing the interactions between human and bot accounts on social media; (vii)
26 characterizing the "Macron Leaks" social media disinformation campaign surrounding the 2017
27 French Presidential election; (viii) characterizing information contagion manipulation by means

Highly Confidential (Competitor)

of social media bots; (ix) detecting early warnings of cyberattacks and cyberthreats on social media; (x) unveiling the use of bots to bolster negative and inflammatory content in the context of the 2017 Catalan referendum; (xi) analyzing the digital traces of political manipulation on social media by the Russian Internet Research Agency (IRA); (xii) developing deep learning based techniques for bot detection; (xiii) measuring spam and the effect of bots on information diffusion on social media; (xiv) unveiling the perils of manipulation of social media discourse surrounding the 2018 US Midterm elections; (xv) characterizing the behavior of Russian IRA trolls operated social media influence campaigns; (xvi) charting the history of digital spam on social media services; (xvii) characterizing linguistic cues to deception used by spammers and trolls on social media; (xviii) comparative analysis of bot partisan behavior on social media; (ixx) characterizing the evolution of bot behavior over election-related events; (xx) charting the landscape of social media cryptocurrency manipulation; (xxi) developing a reinforcement learning framework to detect spam behavior on social media; (xxii) developing a temporal behavioral dynamics model to detect bot behavior on social media; (xxiii) tracking social media related COVID-19 conversation since the inception of the pandemic to track potential manipulation and abuse; (xxiv) characterizing the implications of political polarization surrounding COVID-19 in the US-based social media discussion; (xxv) building a multimodal dataset for information credibility studies of social media and news sources; (xxvi) characterizing social media manipulation by bots and influence campaigns in the context of the 2020 US Presidential Election; (xxvii) unveiling COVID-19 misinformation influencing election-related discourse on social media; (xxviii) characterizing COVID-19 vaccine hesitancy and anti-vaccine campaigns spreading on social media; (xxix) charting the effects of echo chambers in the context of COVID-19 on social media; (xxx) identifying coordinated accounts on social media via hidden influence modeling; (xxxi) characterizing online engagement with election-related disinformation and conspiracies on social media; (xxxii) auditing the effects of algorithmic bias in news-feed ranking on social media; (xxxiii) constructing a large-scale annotated dataset of social media misinformation campaigns; (xxxiv) comparative analysis of bots and humans during the COVID-19 pandemic on social media;

Highly Confidential (Competitor)

(xxxv) surveying influence campaigns during the Stop Asian Hate Movement on social media; (xxxvi) tracking the discourse and influence campaigns on the war between Ukraine and Russia on social media; (xxxvii) charting the landscape of information manipulation in the context of the Israel Hamas conflict; (xxxviii) inferring political leaning of social media users from linguistics and psychographic cues; (xxxix) modeling of user attraction to politically extreme content on social media; (xl) understanding the emotional variance of geolocated social media posts; (xli) studying the effect of user exposure to propaganda on social media; (xlii) modeling the interconnected nature of content moderation and exposure to harmful content on social media; (xliii) unveiling the effects of content censorship by studying Chinese social media; (xliv) modeling the mechanism of adoption and motivators to produce toxic content online; (xlv) modeling the adoption of information after minimal exposure on social media; (xlvi) using AI moderators to improve the quality and limit harmful content exposure online. Since 2015, I have developed and used a plethora of machine learning and AI techniques to study a multitude of problems on social media services, including but not limited to (i) detecting online extremism and radical propaganda by the ISIS terrorist group on social media; (ii) detecting spam in electronic-cigarette advertising campaigns on social media; (iii) unveiling manipulation of social media discourse revolving around the 2016 US Presidential Election; (iv) automatically detecting bots by means of an ensemble of classifiers, known as “Bot Or Not” and later “Botometer” framework; (v) participating in the Defense Advanced Research Projects Agency social media bot challenge and ranking in the top 3 teams for accuracy of bot detection; (vi) characterizing the interactions between human and bot accounts on social media; (vii) characterizing the “Macron Leaks” social media disinformation campaign surrounding the 2017 French Presidential election; (viii) characterizing information contagion manipulation by means of social media bots; (ix) detecting early warnings of cyberattacks and cyberthreats on social media; (x) unveiling the use of bots to bolster negative and inflammatory content in the context of the 2017 Catalan referendum; (xi) analyzing the digital traces of political manipulation on social media by the Russian Internet Research Agency (IRA); (xii) developing deep learning based techniques for bot detection; (xiii)

Highly Confidential (Competitor)

1 measuring spam and the effect of bots on information diffusion on social media; (xiv) unveiling
2 the perils of manipulation of social media discourse surrounding the 2018 US Midterm elections;
3 (xv) characterizing the behavior of Russian IRA trolls operated social media influence campaigns;
4 (xvi) charting the history of digital spam on social media services; (xvii) characterizing linguistic
5 cues to deception used by spammers and trolls on social media; (xviii) comparative analysis of bot
6 partisan behavior on social media; (ixx) characterizing the evolution of bot behavior over election-
7 related events; (xx) charting the landscape of social media cryptocurrency manipulation; (xxi)
8 developing a reinforcement learning framework to detect spam behavior on social media; (xxii)
9 developing a temporal behavioral dynamics model to detect bot behavior on social media; (xxiii)
10 tracking social media related COVID-19 conversation since the inception of the pandemic to track
11 potential manipulation and abuse; (xxiv) characterizing the implications of political polarization
12 surrounding COVID-19 in the US-based social media discussion; (xxv) building a multimodal
13 dataset for information credibility studies of social media and news sources; (xxvi) characterizing
14 social media manipulation by bots and influence campaigns in the context of the 2020 US
15 Presidential Election; (xxvii) unveiling COVID-19 misinformation influencing election-related
16 discourse on social media; (xxviii) characterizing COVID-19 vaccine hesitancy and anti-vaccine
17 campaigns spreading on social media; (xxix) charting the effects of echo chambers in the context
18 of COVID-19 on social media; (xxx) identifying coordinated accounts on social media via hidden
19 influence modeling; (xxxi) characterizing online engagement with election-related disinformation
20 and conspiracies on social media; (xxxii) auditing the effects of algorithmic bias in news-feed
21 ranking on social media; (xxxiii) constructing a large-scale annotated dataset of social media
22 misinformation campaigns; (xxxiv) comparative analysis of bots and humans during the COVID-
23 19 pandemic on social media; (xxxv) surveying influence campaigns during the Stop Asian Hate
24 Movement on social media; (xxxvi) tracking the discourse and influence campaigns on the war
25 between Ukraine and Russia on social media; (xxxvii) charting the landscape of information
26 manipulation in the context of the Israel Hamas conflict; (xxxviii) inferring political leaning of
27 social media users from linguistics and psychographic cues; (xxxix) modeling of user attraction to
28

Highly Confidential (Competitor)

1 politically extreme content on social media; (xl) understanding the emotional variance of
2 geolocated social media posts; (xli) studying the effect of user exposure to propaganda on social
3 media; (xlvi) modeling the interconnected nature of content moderation and exposure to harmful
4 content on social media; (xlvi) unveiling the effects of content censorship by studying Chinese
5 social media; (xlv) modeling the mechanism of adoption and motivators to produce toxic content
6 online; (xlv) modeling the adoption of information after minimal exposure on social media; (xlvii)
7 using AI moderators to improve the quality and limit harmful content exposure online.

8 10. On November 7, 2016, I published the paper titled *Social bots distort the 2016 US*
9 *Presidential election online discussion*. This is the only peer-reviewed research paper to appear
10 before the 2016 US Presidential Election of November 8, 2016. The work shows early signs of
11 foreign state-sponsored operations on social media by means of bots and spam accounts. This work
12 was utilized in the US Senate's investigation into the 2016 election and helped expose the now
13 well-known foreign-sponsored efforts to interfere with the online political discussion.

14 11. My research has informed policy and regulation at both the federal and state level.
15 It has been referenced in laws and bills such as (i) the California "SB 1001 Bots: disclosure" that
16 was ratified into law as the B.O.T. Act ("Bolstering Online Transparency") as of July 2019 and
17 (ii) the "S.2125 - Bot Disclosure and Accountability Act of 2019", introduced by former Senate
18 Judiciary Committee Ranking Member Dianne Feinstein.

19 12. Since 2015, I have published over 200 peer-reviewed research papers. As of June
20 4, 2025, my research has been cited over 25,000 times according to Google Scholar. By this
21 measure of impact, I am one of the top 5 most cited researchers worldwide in the area of "Human-
22 Centered AI," top 20 most cited in "Social Computing" and "Computational Social Science," and
23 top 30 most cited in the areas of "Social Media," and top 200 most cited researchers in the broad
24 area of "Data Science."

25 13. Since 2015, I have secured over \$60 million in research funding, the vast majority
26 of which is funded by the US Department of Defense (DOD), Defense Advanced Research Projects
27 Agency, or government intelligence agencies (e.g., Intelligence Advanced Research Projects

Highly Confidential (Competitor)

1 Activity). This funding has supported research in social media manipulation, bot detection,
2 influence campaign detection, and online service abuse. In addition, certain technologies I have
3 developed have been transitioned, or are in the course of transitioning, into DOD operations by
4 means of Technology Transition awards. Most prominently, the Twitter bot detection techniques
5 I invented have been transitioned to the US government and defense agencies under the Office of
6 Naval Research SBIR/STTR Contract, “DOISAC: Detecting Orchestrated Information and
7 Synthetic Account Campaigns,” of which I was the principal investigator.

8 **C. Publications**

9 14. I have authored, co-authored, and contributed to over 250 academic papers and
10 publications, mostly in data mining and machine learning, particularly with respect to social media
11 networks. A copy of my CV listing my publications and contributions is attached as Exhibit A to
12 this report.

13 **II. Plaintiffs’ Allegations and Summary of Opinions**

14 A. Plaintiffs’ Allegations Regarding Content Moderation

15 15. Plaintiffs allege that Meta has knowingly failed to implement adequate safeguards
16 that prevent, detect, or remove allegedly harmful content on Meta’s family of apps.² Plaintiffs
17 allege that such failure has led to a mental health crisis among Meta’s users, especially teens.³

18 B. Nature of Assignment and Summary of Opinions

19 16. I have been retained by Covington & Burling LLP, counsel for Meta Platforms,
20 Inc., to provide my expert opinions on Meta’s content moderation systems, practices, and related
21 issues on Facebook and Instagram. In particular, I have been asked to assess the reasonableness of
22 Meta’s content moderation systems in light of published academic literature and the challenges of
23 maintaining a content moderation system on social media services at the scale a company like
24 Meta operates. My review of Facebook and Instagram’s content moderation systems is consistent

25

² Master Complaint (Personal Injury) at 46-47, JCCP No. 5255 (Cal. Super. Ct. May 15, 2023).

26 ³ Master Complaint (Personal Injury) *passim*, *Social Media Cases*, JCCP No. 5255 (Cal. Super.
27 Ct. May 15, 2023).

Highly Confidential (Competitor)

1 with how I review such technologies for my research.

2 17. I am being compensated at my customary rate of \$1,200.00 per hour. My
3 compensation is not contingent on the outcome of this litigation.

4 18. I have reviewed the expert reports of Drs. Bagot, Cingel, Christakis, Goldfield,
5 Lembke, Mojtabai, Murray, Telzer and Twenge disclosed on April 18, 2025 and May 16, 2025.
6 These experts have, in general, made high-level allegations about the adequacy of Meta's content
7 moderation efforts, and to the extent that they have done so, the opinions set forth in this report
8 are responsive. I understand that in subsequent rounds of expert reports, Plaintiffs' experts (those
9 above or potentially additional experts) may provide more specific opinions on Meta's content
10 moderation systems, and I reserve the right to respond further to any such opinions.

11 19. I have also evaluated a purported stress test done by Mr. Arturo Bejar. I based my
12 evaluation on my extensive experience in algorithmic auditing, having performed several such
13 studies myself on social media services such as Twitter/X, TikTok, Reddit, and others.

14 20. In summary, I have developed the following opinions based upon my review of the
15 relevant academic literature, filings in this case, documents and testimony produced in discovery,
16 written discovery, and my knowledge and expertise. Attached as Exhibit B, please see a complete
17 list of the sources I have reviewed to develop my opinions. For details of the cases in which I have
18 testified within the last four years, please see Exhibit C.

19 21. Meta has implemented one of the most technically advanced and operationally
20 mature content moderation infrastructures in the industry, in which it has invested substantial
21 resources to develop and continually refine.

22 22. Meta's Community Standards Enforcement Reports ("CSERs") have provided a
23 scientifically sound and statistically robust mechanism for evaluating the effectiveness of content
24 moderation.

25 23. Meta's content moderation systems prioritize the removal of high-severity content,
26 such as suicide & self-injury content and child safety content, and demonstrate exceptionally high
27 proactive detection rates in these categories.

Highly Confidential (Competitor)

24. The company has demonstrated a transparent and iterative approach to measuring, refining, and reporting the performance of its content moderation systems.

25. Meta's content moderation systems are not static but are instead designed to evolve in response to emerging harms, adversarial tactics, cultural considerations, and regulatory constraints.

26. Allegations that Meta's recommendation systems promote harmful content are not grounded in scientifically valid methodologies.

27. Meta’s content moderation systems do not and could not lead to a perfect result in which no policy-violating content is present on Meta’s services—that would be, in my opinion, technically impossible given the immense challenges associated with the massive volume and varied nature of the content on the services, not to mention the efforts of some adversarial, oftentimes criminal, users to try to circumvent Meta’s content moderation systems and intentionally post policy-violating content. But, in my opinion, Meta’s content moderation systems are robust, reasonably designed, and continuously improving.

III. Background

A. Overview of Content Moderation on Social Media Services

28. Social media provides a service for users to create and post content and to communicate with other users. Most social media companies, including Meta, create and enforce policies to prohibit or limit certain types of content on the services. This report refers to the systems and processes used by Meta to enforce such policies as “content moderation.”

29. Content moderation, in the context of social media, is multifaceted and ever changing, involving various stakeholders, including service operators, users, regulators, and society at large. It is an interplay between the technical mechanisms that underpin the services and the policy frameworks that guide their use and governance.

30. Different services define content moderation differently, which at times can make comparing policies or technological solutions across different service providers extremely

Highly Confidential (Competitor)

1 complex.⁴ Irrespective of the definition, content moderation on user generated platforms is an
2 important aspect of the user experience because it dictates the content a user can post and the
3 content a user is likely to see.

4 31. In the initial era of social media services (in the early 2000s), services had relatively
5 small user bases with a relatively limited amount of content being shared. As such, content
6 moderation took place on a relatively small scale and often relied upon simple rule-based
7 approaches.⁵ Typically, small teams of human curators manually reviewed content,⁶ a strategy that
8 was not sustainable in the long term for growing platforms, in part because of sheer scalability
9 challenges, and in part due to transparency and power imbalance issues, with the potential for other
10 issues such as censorship or collusion among users.⁷ Although rudimentary, these early systems
11 established the foundations for what would become more sophisticated moderation systems.⁸

12 32. As user bases grew into the millions, the task of content moderation and feed
13 curation quickly surpassed human review scalability due to the growing volume of content being
14 shared. To address scale, services turned to automated systems as basic machine learning began to
15 take root.⁹ These early systems, however, were often simplistic and usually failed to understand

16 ⁴ See Pierri, Francesco, and Stefano Ceri. *False news on social media: a data-driven survey*.
17 ACM Sigmod Record 48.2 (2019): 18-27; Sharma, Karishma, et al. *Combating fake news: A
18 survey on identification and mitigation techniques*. ACM Transactions on Intelligent Systems
and Technology (TIST) 10.3 (2019): 1-42.

19 ⁵ Jiang, J. A., Middler, S., Brubaker, J. R., & Fiesler, C. (2020, October). *Characterizing
20 community guidelines on social media platforms*. In Companion Publication of the 2020
Conference on Computer Supported Cooperative Work and Social Computing (pp. 287-291).

21 ⁶ See Gillespie, T. (2018). *Custodians of the Internet: Platforms, Content Moderation, and the
22 Hidden Decisions That Shape Social Media*. Yale University Press, which documents the role of
early human moderation on platforms like YouTube and Facebook.

23 ⁷ See Jhaver, S., Bruckman, A., & Gilbert, E. (2019). *Does transparency in moderation really
matter?* CHI Conference on Human Factors in Computing Systems.

24 ⁸ Riedl, M. J. (2020). *Content moderation and volunteer participation*. In The Routledge
25 Encyclopedia of Citizen Media (pp. 93-98). Routledge.

26 ⁹ Gorwa, R., Binns, R., & Katzenbach, C. (2020). Algorithmic content moderation: Technical
27 and political challenges in the automation of platform governance. *Big Data & Society*, 7(1),
2053951719897945.

1 **Highly Confidential (Competitor)**

1 complex social nuances and context, sometimes leading to public controversies around censorship
2 and bias.¹⁰

3 33. In the early 2010s, a surge in content creation brought new challenges to content
4 moderation. Incidents around the abuse of social media services for nefarious purposes were on
5 the rise (e.g., malicious targeting of certain groups by attributing propaganda to members of those
6 groups).¹¹

7 34. Meta invested early in automated defense systems. The Facebook Immune
8 System¹² (“FIS”), launched by early 2011, represents Meta’s first attempt at a real-time integrity
9 infrastructure, designed to detect and mitigate a wide range of abuses, including phishing, spam,
10 and policy-violating content, across billions of daily user actions. Also in 2011, Meta began using
11 PhotoDNA, a technology used to detect known signals of child sexual imagery, on its services.¹³
12 Meta was one of the first large companies to broadly implement the technology and encouraged
13 and trained other companies on how to implement it.¹⁴

14 35. Later, social media companies—Meta chief among them—began to invest in more
15 advanced AI systems capable of detecting nuances in language and imagery, and partnerships with
16 fact-checkers became more common. This era marked the beginning of a more proactive approach
17 to content moderation.

18 **B. Challenges to Content Moderation on Social Media Services**

19 36. As social media companies began to operate globally, their services started to face
20

21 ¹⁰ See Klonick, Kate. *The new governors: The people, rules, and processes governing online*
22 *speech*. Harv. L. Rev. 131 (2017): 1598.

23 ¹¹ See “Disinformation Nation: Social Media’s Role in Promoting Extremism and
24 Misinformation”, 117th Congress (2021-2022), <https://www.congress.gov/event/117th-congress/house-event/111407/text>

25 ¹² See Stein, T., Chen, E., & Mangla, K. (2011, April). *Facebook immune system*. In *Proceedings of the 4th workshop on social network systems*.
26 <https://css.csail.mit.edu/6.858/2014/readings/facebook-immune.pdf>

27 ¹³ META3047MDL-034-00098929.

28 ¹⁴ META3047MDL-034-00098929.

Highly Confidential (Competitor)

1 significant challenges in adapting their policies and corresponding content moderation technology
 2 to shifting global priorities as to what harms are most pressing, as well as to the multijurisdictional
 3 regulatory changes that have followed. Coupled with and exacerbating socio-political challenges
 4 are the technical hurdles that services must overcome in order to respond to public and political
 5 concerns.

6 **1. Socio-political Challenges**

7 37. Societal norms and cultural perspectives shape what is deemed harmful, and these
 8 perspectives shift over time and across regions. This necessitates nuanced understandings of
 9 cultural contexts, which has led to the incorporation of regional expertise in content policy teams
 10 and the development of more localized moderation practices.¹⁵ Certain content that was once
 11 considered benign can become problematic or harmful as social awareness evolves, and regulatory
 12 bodies respond in kind. Content moderation policies must constantly adapt to these shifting
 13 definitions of harm to remain relevant and effective.¹⁶

14 **a. What harm is the most pressing is always evolving**

15 38. Global events, political climates, and public discourse often dictate what perceived
 16 public harm is the most urgent.¹⁷ Meta must maintain agile and responsive policies that
 17 dynamically adjust to these evolving concerns.

18 39. Academic research supports this observation, highlighting the fluid nature of
 19 harmful content definitions and the challenges they pose for content moderation policies.¹⁸ For

20 ¹⁵ See Tworek, Heidi. *History explains why global content moderation cannot work.* (2021).
 21 <https://www.brookings.edu/articles/history-explains-why-global-content-moderation-cannot-work/>

22 ¹⁶ Scheuerman, M. K., Jiang, J. A., Fiesler, C., & Brubaker, J. R. (2021). *A framework of severity*
 23 *for harmful content online.* Proceedings of the ACM on Human-Computer Interaction,
 24 5(CSCW2), 1-33.

25 ¹⁷ Chan, A. J., Redondo García, J. L., Silvestri, F., O'Donnell, C., & Palla, K. (2023). *Enhancing*
 26 *Content Moderation with Culturally-Aware Models.* arXiv e-prints, arXiv-2312.

27 ¹⁸ Arora, A., Nakov, P., Hardalov, M., Sarwar, S. M., Nayak, V., Dinkov, Y., ... & Augenstein, I.
 28 (2023). *Detecting harmful content on online platforms: what platforms need vs. where research*
 29 *efforts go.* ACM Computing Surveys, 56(3), 1-17.

Highly Confidential (Competitor)

1 instance, Dr. Amit Sheth and collaborators authored a prominent study that emphasizes that online
2 problematic content is multi-dimensional and context-sensitive, making its detection
3 challenging.¹⁹ The authors argue that definitions of toxic content extend beyond traditional
4 classifications, necessitating continuous adaptation of moderation policies to address emerging
5 forms of harm. Similarly, Schöpke-Gonzalez and collaborators discussed how varying definitions
6 of harm affect the annotation of data used in content moderation.²⁰ The study reveals that
7 annotators do not use terms like “hateful,” “offensive,” and “toxic” interchangeably, indicating
8 that perceptions of harmful content are subjective and evolving. This underscores the need for
9 moderation policies to be flexible and responsive to changing societal attitudes. Furthermore,
10 Scheuerman and collaborators presented a theoretical framework for understanding the severity of
11 online harms. Their study identifies various types of harm and dimensions along which their
12 severity can be assessed, highlighting the complex and evolving nature of harmful content. The
13 authors suggest that content moderation policies must be adaptable to effectively address the
14 shifting landscape of online harm.

15 40. As another example, research indicated that the COVID-19 pandemic led to a
16 significant increase in anti-Asian hate speech on social media services. A study by researchers at
17 Georgia Tech and Virginia Tech²¹ analyzed over 200 million Tweets and found a notable rise in
18 anti-Asian sentiment during the pandemic. The researchers observed that hate speech often surged
19 following major COVID-19 news events, highlighting the dynamic nature of online hostility in
20 response to global crises. Another study by researchers at CUNY examined the prevalence of anti-

22 ¹⁹ See Sheth, Amit, Valerie L. Shalin, and Ugur Kursuncu. *Defining and detecting toxicity on*
23 *social media: context and knowledge are key.* Neurocomputing 490 (2022): 312-318.

24 ²⁰ Schöpke-Gonzalez, Angela, et al. *How We Define Harm Impacts Data Annotations:*
25 *Explaining How Annotators Distinguish Hateful, Offensive, and Toxic Comments.* arXiv preprint
arXiv:2309.15827 (2023).

26 ²¹ He, Bing, et al. “Racism is a virus: Anti-Asian hate and counter-speech in social media during
the COVID-19 crisis.” Proceedings of the 2021 IEEE/ACM international conference on
advances in social networks analysis and mining. 2021.

Highly Confidential (Competitor)

Asian hate speech on Twitter.²² The findings revealed a significant increase in such content, particularly during the early stages of the pandemic.²³ Other studies further corroborated the escalation of anti-Asian sentiment on social media.²⁴ The pandemic intensified negative attitudes toward Asians, as evidenced by increased usage of racial slurs and derogatory language on services like Twitter. This trend is not limited to public health crises: a study by Olteanu and collaborators examined social media reactions following attacks involving Arabs and Muslims.²⁵ The research observed an increase in online hate speech, particularly messages advocating violence, directed at these communities after such events.

41. These studies underscore the need for effective detection and moderation strategies to address evolving online harms and collectively demonstrate that global events can significantly influence the emergence and prioritization of specific online harms, necessitating adaptive and responsive content moderation policies.

b. Ever-evolving regulatory challenges

42. The regulatory landscape for content moderation is complex and varies across jurisdictions, and this presents additional challenges for online services to perform content moderation. Of the thousands of regulations at the regional, county, and local level, some of the most significant regulatory regimes impacting content moderation efforts are described below:

²² Toliyat, Amir, et al. “Asian hate speech detection on Twitter during COVID-19.” *Frontiers in Artificial Intelligence* 5 (2022): 932381.

²³ To be sure, this content was by no means exclusive to social media platforms; in some circumstances, prominent politicians were fanning anti-Asian sentiments across other media platforms. *See, e.g.*, Itkowitz, C. (2020, Jan. 23). *Trump again uses racially insensitive term to describe coronavirus*. Washington Post. https://www.washingtonpost.com/politics/trump-again-uses-kung-flu-to-describe-coronavirus/2020/06/23/0ab5a8d8-b5a9-11ea-aca5-ebb63d27e1ff_story.html.

²⁴ Lu, Runjing, and Sophie Yanying Sheng. “How racial animus forms and spreads: Evidence from the coronavirus pandemic.” *Journal of Economic Behavior & Organization* 200 (2022): 82-98.

²⁵ Olteanu, Alexandra, et al. “The effect of extremist violence on hateful speech online.” *Proceedings of the international AAAI conference on web and social media*. Vol. 12. No. 1. 2018.

1 **Highly Confidential (Competitor)**

2 43. **General Data Protection Regulation (“GDPR”):** The EU’s GDPR imposes strict
3 guidelines on the processing of personal data, emphasizing individual privacy rights. For content
4 moderation, this means that services must carefully manage user data, particularly concerning data
5 retention for AI model training.²⁶ The GDPR’s principles of data minimization and purpose
6 limitation can limit the extent to which user data is stored and utilized, potentially hindering further
7 development and refinement of AI moderation tools.²⁷

8 44. **Digital Services Act (“DSA”):** The EU’s 2022 DSA introduces comprehensive
9 obligations for online services, including enhanced transparency requirements and stricter
10 accountability for algorithmic content recommendations. Services are mandated to provide
11 detailed explanations for content removal decisions and help ensure that their algorithms do not
12 contribute to the dissemination of illegal content. This necessitates significant adjustments in
content moderation practices to comply with the new standards.²⁸

13 45. **Children’s Online Privacy Protection Act (“COPPA”):** The US’s COPPA
14 prohibits in certain cases the collection and retention of “personal information” from children
15 under the age of thirteen without verifiable parental notice and consent.²⁹

16 46. The scope of these and other regulatory regimes places significant hurdles before
17 any social media service that moderates content, including Meta. Some of the most significant
18 challenges impacting content moderation efforts are described below:

19 47. **Limits to Human Review:** Privacy regulations may restrict the use of human
20 reviewers in content moderation due to concerns over data protection. For instance, the GDPR

21

²⁶ *The impact of the GDPR on artificial intelligence.* (n.d.). Securiti.ai. Retrieved June 10, 2025,
22 from <https://securiti.ai/impact-of-the-gdpr-on-artificial-intelligence>.

23 ²⁷ *The impact of the GDPR on artificial intelligence.* (n.d.). Securiti.ai. Retrieved June 10, 2025,
24 from <https://securiti.ai/impact-of-the-gdpr-on-artificial-intelligence>.

25 ²⁸ Nunziato, D. C. (2023). The Digital Services Act and the Brussels Effect on Platform Content
26 Moderation. Chi. J. Int’l L., 24, 115.

27 ²⁹ Harmonizing laws to improve enforcement of human rights crimes. (n.d.). Thomson Reuters.
28 Retrieved June 10, 2025, from <https://www.thomsonreuters.com/en-us/posts/human-rights-crimes/harmonizing-laws/>.

Highly Confidential (Competitor)

1 mandates that data processing involving human intervention must uphold user privacy rights,
 2 potentially limiting the extent to which human moderators can access and assess user content. This
 3 can challenge services' abilities to effectively manage nuanced content that AI systems may
 4 struggle to interpret.³⁰

5 **48. *Restrictions on AI-Enhanced Moderation Tools:*** Privacy laws like the GDPR may
 6 limit the ability to use AI-enhanced moderation tools that rely on extensive user data for training
 7 and improvement. The requirement for explicit user consent and the right to data deletion can
 8 reduce the data available for AI models, potentially diminishing their effectiveness in detecting
 9 harmful content.³¹ In Europe, the Privacy Directive also prohibits processing of private messages
 10 without GDPR-level consent, also diminishing the data available for AI content moderation
 11 models.³² Likewise, in the US, Meta has faced significant regulatory hurdles when developing AI
 12 technology that can reliably identify users under the age of thirteen at scale,³³ and COPPA's
 13 barriers to the collection and use of personal information from under-thirteen users have added to
 14 those challenges.³⁴

15 **49. *Balancing Transparency and Security:*** Transparency mandates, such as those in
 16 the DSA, require services to disclose information about their content moderation algorithms.
 17 While this promotes accountability, it also poses risks, as overly detailed disclosures could enable

18 ³⁰ Information Commissioner's Office. (n.d.). What is content moderation and how does it use
 19 personal information?. Retrieved June 10, 2025, from <https://ico.org.uk/for-organisations/uk-gdpr-guidance-and-resources/online-safety-and-data-protection/content-moderation-and-data-protection/what-is-content-moderation-and-how-does-it-use-personal-information/>.

21 ³¹ See Masnick, M. (2018, January 30). Unintended consequences of EU's new internet privacy
 22 rules: Facebook won't use AI to catch suicidal users. Techdirt.
<https://www.techdirt.com/2018/01/30/unintended-consequences-eus-new-internet-privacy-rules-facebook-wont-use-ai-to-catch-suicidal-users/>.

23 ³² See European Data Protection Board. (2020, May 4). *Guidelines 05/2020 on consent under
 24 Regulation 2016/679*.
https://www.edpb.europa.eu/sites/default/files/files/file1/edpb_guidelines_202005_consent_en.pdf; GDPR: Consent. (n.d.). Intersoft Consulting. Retrieved June 10, 2025, from <https://gdpr-info.eu/issues/consent/>.

27 ³³ META3047MDL-065-00042982.

28 ³⁴ META3047MDL-034-00027362.

Highly Confidential (Competitor)

malicious actors to game the system, circumventing moderation efforts. Services must carefully balance the need for transparency with the imperative to maintain the integrity of their moderation systems.

50. ***Tension Between Free Speech and Content Removal:*** Regulations often require the removal of illegal or harmful content, but this can conflict with free speech protections, particularly in jurisdictions like the United States, where the First Amendment provides robust speech rights. Services must navigate these conflicting expectations, ensuring compliance with content removal mandates while respecting users' rights to free expression.³⁵

51. ***Conflicting Government Expectations:*** Different governments impose varying expectations on content moderation. The European Union, through regulations like the DSA, demands stricter moderation to prevent harm, whereas in the United States, there is often advocacy for less intervention to uphold free speech. This creates a complex regulatory environment for services operating globally, as they must tailor their moderation practices to comply with diverse legal frameworks.¹⁵

2. Technical Challenges for State of the Art Artificial Intelligence Systems

52. From a technical standpoint, the accuracy and effectiveness of content moderation systems in detecting and filtering out harmful content are ongoing subjects of research. Despite advances in natural language processing and image recognition technologies, algorithms still struggle with understanding context, sarcasm, and subtlety, leading to both over-censorship and underenforcement.³⁶

³⁵ For example, in Thailand, laws exist that prohibit insulting the monarchy, which stands in stark contrast to US First Amendment protections associated with criticism of the government. BBC News. (2017, October 6). *Lese-majeste explained: How Thailand forbids insult of its royalty*. <https://www.bbc.com/news/world-asia-29628191>.

³⁶ See Gorwa, Robert, Reuben Binns, and Christian Katzenbach. *Algorithmic content moderation: Technical and political challenges in the automation of platform governance*. *Big Data & Society* 7.1 (2020): 2053951719897945.

Highly Confidential (Competitor)

1 53. To be sure, artificial intelligence systems for text classification have advanced
2 substantially in recent years, largely due to the emergence and adoption of transformer-based
3 architectures. Introduced by Vaswani et al. in 2017,³⁷ the transformer model revolutionized natural
4 language processing (NLP) by replacing traditional recurrent and convolutional networks
5 architectures³⁸ with a self-attention mechanism that allows the model to weigh the relevance of
6 each word in a sequence, regardless of its position. This design enables transformers to better
7 capture long-range dependencies and contextual relationships within text.

8 54. Transformers form the backbone of modern large language models (LLMs), such
9 as BERT (Bidirectional Encoder Representations from Transformers),³⁹ GPT (Generative Pre-
10 trained Transformer),⁴⁰ and RoBERTa (a robustly optimized BERT variant),⁴¹ which have set new
11 benchmarks across a wide array of NLP tasks, including sentiment analysis, hate speech detection,
12 and content moderation. These models are pre-trained on vast corpora of text using self-supervised
13 objectives, and then fine-tuned on specific downstream tasks using labeled data. Their
14 bidirectional (or autoregressive) architectures allow them to capture both left and right context,
15 improving semantic understanding and disambiguation of meaning.

16 55. In the context of content moderation, transformer-based LLMs excel at parsing
17 subtle linguistic signals such as sarcasm, idioms, code-switching, and culturally-specific
18

19
20 ³⁷ Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). *Attention is all you need*. Advances in neural information processing systems, 30.

21 ³⁸ META3047MDL-208-00052347.

22 ³⁹ Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019, June). *BERT: Pre-training of deep*
23 *bidirectional transformers for language understanding*. In Proceedings of the 2019 conference of
24 the North American chapter of the association for computational linguistics: human language
technologies, volume 1 (long and short papers) (pp. 4171-4186).

25 ⁴⁰ Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). *Improving language*
understanding by generative pre-training.

26 ⁴¹ Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019). *RoBERTa:*
A robustly optimized BERT pretraining approach. arXiv preprint arXiv:1907.11692.

Highly Confidential (Competitor)

1 references that traditional machine learning models often miss.⁴² Meta scientists adapted
2 transformer-based models to their content moderation challenges, such as Meta's multilingual
3 model like XLM-R,⁴³ to manage content across diverse languages and dialects. Furthermore,
4 research by Meta's AI teams has led to the development of more efficient variants, such as the
5 Linformer, which reduces the computational overhead associated with processing long sequences
6 of text, making real-time moderation more scalable.⁴⁴

7 56. Nevertheless, detecting and moderating harmful content on social media services
8 like Meta's remains a multifaceted technical challenge. The main technical challenges associated
9 with AI automated content moderation systems include the complexity of language, the scale and
10 volume of content generated daily, resource intensiveness, identifying false positives and
11 negatives, evolving adversarial threats (e.g., deepfakes), and the growing generation and alteration
12 of existing content using AI.⁴⁵

13 57. The technical challenges of detecting and moderating harmful content on social
14 media services are continuously changing. Advances in AI, including those spearheaded by Meta's
15
16
17

18⁴² Chan, A. J., Redondo García, J. L., Silvestri, F., O'Donnell, C., & Palla, K. (2023). *Enhancing
Content Moderation with Culturally-Aware Models*. arXiv e-prints, arXiv-2312.

19⁴³ Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., ... &
20 Stoyanov, V. (2019). Unsupervised cross-lingual representation learning at scale. arXiv preprint
arXiv:1911.02116.

21⁴⁴ Wang, S., Li, B. Z., Khabsa, M., Fang, H., & Ma, H. (2020). Linformer: Self-attention with
22 linear complexity. arXiv preprint arXiv:2006.04768.

23⁴⁵ Gorwa, R., Binns, R., & Katzenbach, C. (2020). "Algorithmic content moderation: Technical
24 and political challenges in the automation of platform governance." *Big Data & Society.*; Udupa,
S., Maronikolakis, A., & Wisiorek, A. (2023). Ethical scaling for content moderation: Extreme
25 speech and the (in) significance of artificial intelligence. *Big Data & Society*, 10(1),
20539517231172424; Huertas-García, Á., Martín, A., Huertas-Tato, J., & Camacho, D. (2023).
Countering malicious content moderation evasion in online social networks: Simulation and
26 detection of word camouflage. *Applied Soft Computing*, 145, 110552; Matias, J. Nathan. (2019).
"The civic labor of volunteer moderators: Moderation, productivity, and the limits of automation
27 at scale." *New Media & Society*, 21(6), 1271–1287.

1 **Highly Confidential (Competitor)**

1 AI scientists—such as efficient Transformer architectures and multimodal learning⁴⁶—have
2 significantly improved the ability to detect harmful content.

3 **58. Cultural Challenges and Complexity of Language:** Cultural challenges arise when
4 content moderation algorithms fail to recognize contextual nuances, leading to the suppression of
5 culturally significant content or the inadvertent promotion of culturally insensitive material.⁴⁷ It is
6 helpful for services to engage with experts from around the globe and incorporate a diverse set of
7 perspectives to help ensure that algorithms are sensitive to cultural variations and do not enforce a
8 one-size-fits-all approach to content moderation.⁴⁸ Automated systems must navigate the
9 intricacies of human language, including sarcasm, slang, and cultural references, which vary across
10 different languages and dialects.⁴⁹ This complexity makes it challenging for AI models to
11 accurately interpret and classify content.⁵⁰ A recent study published by Matias and colleagues in
12 *New Media & Society* highlights the severe limitations of existing machine learning content
13 moderation methods in understanding such linguistic nuances: They found that classifiers often
14 misinterpreted sarcasm, coded language, and in-group slang, flagging benign speech while failing
15 to detect covert abuse. The study highlighted that models trained without cultural and contextual
16 nuance performed poorly across diverse communities,⁵¹ which can in turn lead to ethical concerns

17
18 ⁴⁶ See Wang, Sinong, Belinda Z. Li, Madien Khabsa, Han Fang, and Hao Ma. “Linformer: Self-
19 attention with linear complexity.” *arXiv preprint arXiv:2006.04768* (2020). Note that this paper
20 has been cited over a thousand times by researchers in the field of AI, substantiating its
21 contribution and importance.

22
23 ⁴⁷ See Gillespie, Tarleton. (2018.). *Custodians of the Internet: Platforms, content moderation,*
24 *and the hidden decisions that shape social media.* Yale University Press

25
26 ⁴⁸ META3047MDL-003-00001478.

27
28 ⁴⁹ Indeed, the same holds true for human reviewers, as differences in cultural norms and
experience may have a negative impact on the accuracy of human review.

24
25 ⁵⁰ Udupa, S., Maronikolakis, A., & Wisiorek, A. (2023). Ethical scaling for content moderation:
26 Extreme speech and the (in) significance of artificial intelligence. *Big Data & Society*, 10(1),
27 20539517231172424.

28
29 ⁵¹ Matias, J. Nathan. (2019). “The civic labor of volunteer moderators: Moderation, productivity,
and the limits of automation at scale.” *New Media & Society*, 21(6), 1271–1287.

Highly Confidential (Competitor)

1 about the fairness of such systems.⁵²

2 59. Furthermore, there are subtleties of meaning in context, related to cultural or in-
 3 group communication norms that pose significant content moderation challenges: often, a word or
 4 image can assume different meanings in different communities. For instance, a picture of train
 5 tracks may have to do with travel (and therefore be innocuous content), or it may relate to SSI.⁵³
 6 While humans who appreciate this context may more quickly tell the difference between these
 7 images, training AI to do so is challenging, particularly at scale. This is even harder on services
 8 like Instagram that are mostly image-based, where content often lacks any textual cues.⁵⁴

9 60. For example, detecting hate speech and other harmful content requires a deep
 10 understanding of context, idioms, and cultural nuances, which even state-of-the-art models
 11 struggle to achieve consistently.⁵⁵ While attention weights offer some insight into model behavior,
 12 decisions made by these models are often opaque, complicating accountability and appeal
 13 processes in automated content enforcement.⁵⁶

14 61. ***Multimodal and Multilingual Content Learning & Integration:*** Harmful content
 15 is not confined to text alone – it can also appear in images, videos, audios, and a combination of
 16 these modalities. Additionally, it can be expressed in various languages⁵⁷ and dialects.⁵⁸ AI

17 ⁵² Udupa, S., Maronikolakis, A., & Wisiorek, A. (2023). Ethical scaling for content moderation:
 18 Extreme speech and the (in) significance of artificial intelligence. *Big Data & Society*, 10(1),
 20539517231172424.

19 ⁵³ META3047MDL-001-00000112.

20 ⁵⁴ META3047MDL-001-00000112.

21 ⁵⁵ Meta. (n.d.). *AI advances to better detect hate speech*. Retrieved June 10, 2025,
 22 from <https://ai.meta.com/blog/ai-advances-to-better-detect-hate-speech/>.

23 ⁵⁶ Binns, R., Van Kleek, M., Veale, M., Lyngs, U., Zhao, J., & Shadbolt, N. (2018, April). ‘It’s
 24 Reducing a Human Being to a Percentage’ Perceptions of Justice in Algorithmic Decisions. In
 Proceedings of the 2018 Chi conference on human factors in computing systems (pp. 1-14).

25 ⁵⁷ Haider, S., Luceri, L., Deb, A., Badawy, A., Peng, N., & Ferrara, E. (2023, April). Detecting
 26 social media manipulation in low-resource languages. In Companion Proceedings of the ACM
 27 Web Conference 2023 (pp. 1358-1364).

28 ⁵⁸ Nicholas, G., & Bhatia, A. (2023). Toward better automated content moderation in low-
 29 resource languages. *Journal of Online Trust and Safety*, 2(1).

Highly Confidential (Competitor)

systems must integrate data from multiple modalities to accurately detect content that combines text and imagery, such as memes that use visual and textual elements to convey harmful messages.⁵⁹ Combining insights from different data modalities to form a coherent understanding of the content is a complex task that requires advanced algorithms and extensive training data.⁶⁰ Models must also be sensitive to the cultural contexts in which different languages are used to avoid misclassifying content due to cultural misunderstandings.⁶¹ Costs to annotate multimodal data can rise rapidly at the scale that Meta operates.

62. **Resource Intensiveness:** Transformer models are notoriously resource-intensive, requiring substantial computational resources and memory, which can limit their scalability. Costs to deploy and maintain these models can also rise rapidly at the scale Meta operates.⁶² For example, Meta relies upon over 1,000 machine learning models at any one time, to continuously inspect content posted on Instagram.⁶³

63. **Scale and Volume:** Services like Facebook and Instagram process a colossal

⁵⁹ Sharma, S., Alam, F., Akhtar, M. S., Dimitrov, D., Martino, G. D. S., Firooz, H., ... & Chakraborty, T. (2022). Detecting and understanding harmful memes: A survey. arXiv preprint arXiv:2205.04274.

⁶⁰ Lippe, P., Holla, N., Chandra, S., Rajamanickam, S., Antoniou, G., Shutova, E., & Yannakoudakis, H. (2020). A multimodal framework for the detection of hateful memes. arXiv preprint arXiv:2012.12871.

⁶¹ Bui, M. D., von der Wense, K., & Lauscher, A. (2024). Multi3Hate: Multimodal, Multilingual, and Multicultural Hate Speech Detection with Vision-Language Models. arXiv preprint arXiv:2411.03888.

⁶² See Strubell, E., Ganesh, A., & McCallum, A. (2019). “Energy and policy considerations for deep learning in NLP.” ACL. This study quantifies the high energy costs and carbon footprint of training transformer models, such as BERT and GPT-2, highlighting the practical trade-offs of deploying these models at large scale. State of the art models, such as OpenAI GPT4 and next-gen iterations are likely to have orders of magnitude higher costs. Fanatical Futurist. (2025, May). *OpenAI GPT-5 is costing \$500 million per training run and still failing?* <https://www.fanaticalfuturist.com/2025/05/openai-gpt-5-is-costing-500-million-per-training-run-and-still-failing/>.

⁶³ Meta. (2025, May 21). *Journey to 1000 models: Scaling Instagram’s recommendation system.* <https://engineering.fb.com/2025/05/21/production-engineering/journey-to-1000-models-scaling-instagram-recommendation-system/>.

1 **Highly Confidential (Competitor)**

2 volume content daily,⁶⁴ requiring efficient algorithms and robust infrastructure for real-time data
3 management: Meta's AI systems rely on hundreds of different statistical, algorithmic, and
4 machine-learning models at any one time to process, filter, and recommend content.⁶⁵ The
5 immense scale of data and the relentlessness of violations make AI approaches desirable, even
6 inevitable, for content moderation. However, achieving effective moderation at this scale remains
7 a significant challenge.⁶⁶ Despite algorithmic content moderation, Meta receives and processes
8 millions of user reports each week.⁶⁷ The vast amount of content generated on social media
9 services necessitates automated solutions for effective monitoring and management. AI systems
10 must be highly accurate and efficient to operate in real-time across global services and varied data
11 modalities (text, images, videos, audio). This requires continuous improvements in both hardware
12 and software to maintain performance at scale.

13 64. **Data Scarcity for Training:** Building effective AI models for content moderation
14 requires extensive labeled data representing both harmful and benign content. Supervised learning
15 (i.e., learning from examples, the dominant paradigm in AI moderation) relies on such annotations
16 to train models that can generalize effectively to new inputs. However, certain types of violations,
17 such as terrorism-related content, CSAM, or self-harm imagery, are relatively rare and legally
18 restricted, making large-scale collection and labeling infeasible. Some scholars have highlighted
19 how data imbalance and the scarcity of high-quality labeled examples limit the effectiveness of
20

21

⁶⁴ Smith, C. (2013, September 18). *Facebook users are uploading 350 million new photos each day*. Business Insider. <https://www.businessinsider.com/facebook-350-million-photos-each-day-2013-9>.

22
23 ⁶⁵ Meta. (2025, May 21). *Journey to 1000 models: Scaling Instagram's recommendation system*.
24 <https://engineering.fb.com/2025/05/21/production-engineering/journey-to-1000-models-scaling-instagram-recommendation-system/>.

25 ⁶⁶ Gorwa, Robert, Reuben Binns, and Christian Katzenbach. *Algorithmic content moderation: Technical and political challenges in the automation of platform governance*. *Big Data & Society* 7.1 (2020): 2053951719897945.

26
27 ⁶⁷ META3047MDL-003-00144612.

Highly Confidential (Competitor)

hate speech detection models.⁶⁸ Others similarly noted that content moderation models often fail to generalize across platforms or contexts due to lack of representative data, especially for low-resource languages and culturally specific slurs or euphemisms.⁶⁹ Platform policies, privacy law (e.g., GDPR), and ethical limitations introduce major constraints on the collection of training data involving vulnerable populations, such as children or marginalized users.⁷⁰ These constraints prevent the publication or open sharing of many datasets, limiting opportunities for cross-institutional model development or reproducibility. A useful case study is the Facebook Hateful Memes Challenge dataset (Kiela et al., 2020),⁷¹ which combined image-text pairs to model subtle multimodal hate. Despite its innovation, the authors explicitly acknowledged limitations due to annotation complexity, cultural context dependence, and insufficient volume for broader generalization. This underscores how even large, well-designed datasets fall short of solving the scarcity challenge in real-world deployment. Another study on leveraging large-scale multimedia datasets highlighted the difficulty in creating accurate models due to the lack of adequate task-specific training data.⁷² The authors demonstrated that even large, widely-used annotated datasets (such as the Twitter hate speech corpora frequently used to train content moderation systems,

⁶⁸ Schmidt, A., & Wiegand, M. (2017, April). A survey on hate speech detection using natural language processing. In Proceedings of the fifth international workshop on natural language processing for social media (pp. 1-10).

⁶⁹ Vidgen, B., & Derczynski, L. (2020). Directions in abusive language training data, a systematic review: Garbage in, garbage out. Plos one, 15(12), e0243300.

⁷⁰ Gorwa, R., Binns, R., & Katzenbach, C. (2020). Algorithmic content moderation: Technical and political challenges in the automation of platform governance. Big Data & Society, 7(1), 2053951719897945.

⁷¹ Kiela, D., Firooz, H., Mohan, A., Goswami, V., Singh, A., Ringshia, P., & Testuggine, D. (2020). The hateful memes challenge: Detecting hate speech in multimodal memes. Advances in neural information processing systems, 33, 2611-2624.

⁷² Sarris, I., Koutlis, C., Papadopoulou, O., & Papadopoulos, S. (2022, December). Leveraging large-scale multimedia datasets to refine content moderation models. In 2022 IEEE Eighth International Conference on Multimedia Big Data (BigMM) (pp. 125-132). IEEE.

Highly Confidential (Competitor)

1 including DWMW17,⁷³ and FDCL18,⁷⁴ etc.) suffer from systemic biases in how crowd workers
2 label dialect-associated content, particularly African American English (“AAE”). These biases
3 propagate into machine learning models trained on such data, leading to disproportionately high
4 false positive rates for minority dialects. The study emphasizes that task-specific data is not only
5 scarce but often misrepresentative, especially for marginalized users, and calls for more careful
6 annotation protocols, dialect-awareness, and transparency in dataset design to ensure equity in
7 content moderation.

8 65. **Data Quality:** Ensuring the quality and representativeness of training data is
9 essential to avoid biases and improve model performance. AI models can inherit biases present in
10 training data, leading to unfair or discriminatory content moderation decisions. Biased algorithms
11 can disproportionately affect certain demographics, reinforcing existing prejudices: for example,
12 a prominent study empirically demonstrated how racial bias can be introduced and perpetuated in
13 machine learning systems trained for content moderation.⁷⁵ The authors examined two widely-
14 used Twitter datasets annotated for toxic language and uncovered a statistically significant
15 correlations between features of AAE and labels for “offensive” or “abusive” content, finding that
16 annotators (typically drawn from predominantly white crowd worker pools) were more likely to
17 label Tweets written in AAE as toxic, even when those Tweets were benign or in-group vernacular.
18 Left unaddressed, such biases can compromise fairness toward users. Biases can be mitigated by
19 having humans working to ensure representativeness in training data, using dialect-primed or
20

21 ⁷³ Davidson, T., Warmsley, D., Macy, M., & Weber, I. (2017, May). Automated hate speech
22 detection and the problem of offensive language. In Proceedings of the international AAAI
23 conference on web and social media (Vol. 11, No. 1, pp. 512-515).

24 ⁷⁴ Founta, A., Djouvas, C., Chatzakou, D., Leontiadis, I., Blackburn, J., Stringhini, G., ... &
25 Kourtellis, N. (2018, June). Large scale crowdsourcing and characterization of twitter abusive
behavior. In Proceedings of the international AAAI conference on web and social media (Vol.
12, No. 1).

26 ⁷⁵ Sap, M., Card, D., Gabriel, S., Choi, Y., & Smith, N. A. (2019, July). The risk of racial bias in
27 hate speech detection. In Proceedings of the 57th annual meeting of the association for
computational linguistics (pp. 1668-1678).

1 **Highly Confidential (Competitor)**

1 demographically diverse annotation protocols, and auditing model outcomes across identity
2 groups to prevent these harms. For example, Meta has made strides to collect demographic data
3 with an aim toward improving the fairness of its technologies.⁷⁶

4 66. **Generalization:** Models must generalize from training data to real-world scenarios,
5 which can vary significantly in terms of content and context. A recent study discusses the
6 challenges in fine-tuning models to handle diverse and evolving content effectively.⁷⁷ The study
7 explores these challenges using LLMs that are trained not just to label harmful content, but to
8 explain why the content is harmful in plain language. This reasoning-based training helps the
9 models better understand the meaning behind the words, which makes them more accurate and
10 fair, even on content they have not seen before.

11 67. **False Positives and False Negatives:** Balancing sensitivity and specificity in
12 content moderation is particularly challenging because AI models must minimize false positives
13 (i.e., incorrectly flagging benign content) and false negatives (i.e., failing to detect harmful
14 content) among a sea of fluctuating contexts that make up human interaction. Continuous learning
15 and model updates are necessary to improve detection accuracy. A report by *New America*
16 emphasizes the complexity of automated tools in content moderation, noting the challenges that
17 AI researchers face when they have to construct comprehensive datasets that account for the vast
18 fluidity and variances in human language, leading to potential false positives and negatives.⁷⁸ For
19 example, a recent study highlighted how two LLM-based models for content moderation based,
20

21 ⁷⁶ Meta. (n.d.). *Assessing fairness of our products while protecting people's privacy*. Retrieved
22 June 10, 2025, from <https://ai.meta.com/blog/assessing-fairness-of-our-products-while-protecting-peoples-privacy/>.

23 ⁷⁷ Ma, H., Zhang, C., Fu, H., Zhao, P., & Wu, B. (2023). Adapting large language models for
24 content moderation: Pitfalls in data engineering and supervised fine-tuning. arXiv preprint
arXiv:2310.03400.

25 ⁷⁸ New America. (n.d.). *The limitations of automated tools in content moderation*. Retrieved June
26 10, 2025, from <https://www.newamerica.org/oti/reports/everything-moderation-analysis-how-internet-platforms-are-using-artificial-intelligence-moderate-user-generated-content/the-limitations-of-automated-tools-in-content-moderation>.

1 **Highly Confidential (Competitor)**

2 ToxiGen-RoBERTa and RoBERTa-Toxicity-Classifier, based on work done by researchers at
3 Google and Microsoft Research, which are widely adopted, exhibit different kinds of limitations
that can affect false positive and negative results when classifying toxic content:⁷⁹

- 4 • **Arbitrariness of model decisions:** The paper highlights that approximately 30%
5 of English statements might be flagged as "toxic" (a potential false positive if it's
6 legitimate content) by one version of a model, while another equally accurate
7 version of the model might deem it acceptable.
- 8 • **Disparate impact:** The finding that arbitrary moderation decisions are unequally
9 distributed across different demographic groups (e.g., anti-LGBTQ speech
10 receiving more arbitrary decisions) implies that certain groups might be
11 disproportionately affected by both false positives (legitimate speech being
12 removed) and false negatives (harmful content targeting them being allowed).
- 13 • **Limitations in replicating human judgment:** The paper states that models can
14 exhibit high levels of disagreement even on content where human annotators
15 unanimously agree on its toxicity. This suggests that the models are either
16 generating false negatives (missing clearly toxic content) or false positives
17 (flagging content as toxic that humans agree is not), or both, at an inconsistent
18 rate.

19 68. **Evolving Digital Threats:** Malicious actors employ various strategies to evade
20 detection. Strategies such as using deliberate misspellings or code words,⁸⁰ and altered images⁸¹

21

⁷⁹ Gomez, J. F., Machado, C., Paes, L. M., & Calmon, F. (2024, June). Algorithmic arbitrariness
22 in content moderation. In Proceedings of the 2024 ACM Conference on Fairness, Accountability,
23 and Transparency (pp. 2234-2253).

24 ⁸⁰ Wang, W., Huang, J. T., Wu, W., Zhang, J., Huang, Y., Li, S., ... & Lyu, M. R. (2023, May).
25 Mttm: Metamorphic testing for textual content moderation software. In 2023 IEEE/ACM 45th
26 International Conference on Software Engineering (ICSE) (pp. 2387-2399). IEEE.

27 ⁸¹ Jiang, Z., Zhang, J., & Gong, N. Z. (2023, November). Evading watermark-based detection of
28 ai-generated content. In Proceedings of the 2023 ACM SIGSAC Conference on Computer and
Communications Security (pp. 1168-1181).

Highly Confidential (Competitor)

1 are commonly used to bypass AI filters. Spammers, propagandists, and other bad-faith users
2 continually evolve their tactics to circumvent algorithmic detection, necessitating constant updates
3 and improvements to integrity systems.⁸² The adaptability of malicious actors poses a significant
4 technical challenge.

5 69. For example, in recent years, Meta removed a network of malicious actors
6 operating from Russia that relied on large networks of fake accounts that amplified each others'
7 deceptive content.⁸³ Meta must constantly work to identify and stop coordinated efforts to bypass
8 its content moderation systems.⁸⁴

9 70. Adversarial actors often find ways to bypass integrity safeguards, employing tactics
10 such as sophisticated bots, and coordinated inauthentic behaviors (e.g., state-sponsored influence
11 operations, coordinated spam or fraud campaigns). A recent study analyzed how malicious users
12 evade AI-based moderation on platforms like Facebook:⁸⁵ Drawing on a dataset of 500 flagged
13 posts, the authors identified ten evasion tactics (including image alteration, language obfuscation,
14 and contextual manipulation) that exploit weaknesses in automated moderation systems. The study
15 concluded that while AI enables scalable moderation, it remains vulnerable to adversarial tactics
16 and must be paired with human oversight, emphasizing the need for adaptive, context-aware

17 82 See a) Ferrara, Emilio, et al. "The rise of social bots." *Communications of the ACM* 59.7
18 (2016): 96-104; b) Ferrara, Emilio. "The history of digital spam." *Communications of the ACM*
19 62.8 (2019): 82-91; c) Ferrara, Emilio. "GenAI against humanity: Nefarious applications of
20 generative artificial intelligence and large language models." *Journal of Computational Social
Science* (2024): 1-21.1 and d) Huertas-García, Á., Martín, A., Huertas-Tato, J., & Camacho, D.
21 (2023). Countering malicious content moderation evasion in online social networks: Simulation
22 and detection of word camouflage. *Applied Soft Computing*, 145, 110552.

23 83 Meta. (2022, February). *January 2022 coordinated inauthentic behavior report*. Retrieved
24 June 10, 2025, from <https://about.fb.com/news/2022/02/january-2022-coordinated-inauthentic-behavior-report/>.

25 84 E.g., Meta. (2022, February). *January 2022 coordinated inauthentic behavior report*.
26 Retrieved June 10, 2025, from <https://about.fb.com/news/2022/02/january-2022-coordinated-inauthentic-behavior-report/>.

27 85 Malec, L., & Lešetický, J. (2024). *Social media content moderation, censorship and AI
detection evasion techniques*. IDIMT-2024: Changes to ICT, Management, and Business
Processes through AI.

1 **Highly Confidential (Competitor)**

systems to respond effectively to evolving digital threats.

2 **3. Limitations of human reviewers**

3 **71. Diverse Cultural and Linguistic Backgrounds:** One of the main challenges to
4 robust and accurate human review on a global social media service is identifying and recruiting
5 qualified human reviewers across the world who represent a spectrum of diverse cultures and
6 languages. For example, Meta's services are used globally, requiring reviewers who understand
7 local contexts, cultural nuances, and languages to accurately assess content.⁸⁶ What may be
8 considered offensive or harmful in one culture might be acceptable in another. Ensuring that
9 reviewers have this cultural sensitivity is crucial for fair and effective moderation.

10 **72. Geographical Distribution:** A global social media service like Meta requires
11 reviewers in different time zones to help ensure 24/7 coverage. This requires establishing review
12 teams across various regions, which adds to the complexity of recruitment and management.

13 **73. Scale:** Scaling up human review teams involves logistical complexities, such as
14 recruiting, training, and managing a large and diverse workforce made up of around 35,000 content
15 reviewers worldwide.⁸⁷ This requires significant investment in infrastructure and resources, and
16 scaling up pools of reviewers for a particular market can sometimes take several months, if not
17 longer (e.g., sub-Saharan African markets have proven difficult to scale).

18 **74. Accuracy vs. Overenforcement:** Striking a balance between accurately detecting
19 harmful content and avoiding the suppression of legitimate speech is crucial, as overly aggressive
20 filters can stifle free expression and counterspeech.⁸⁸ Indeed, false positives of purportedly illegal
21 content could have serious legal ramifications for wrongly accused users. In this regard, Meta's

22
23 ⁸⁶ See Meta. (n.d.). *Community Standards*. Retrieved June 10, 2025,
24 from <https://transparency.meta.com/policies/community-standards/>; Meta. (n.d.). *How review*
25 *teams work*. Retrieved June 10, 2025,
from <https://transparency.meta.com/enforcement/detecting-violations/how-review-teams-work/>.

26 ⁸⁷ See META3047MDL-003-00144612.

27 ⁸⁸ Meta. (n.d.). *Community Standards*. Retrieved June 10, 2025,
from <https://transparency.meta.com/policies/community-standards/>.

1 **Highly Confidential (Competitor)**

1 posture is to strike a balance between restricting harmful content and promoting expression.⁸⁹

2 * * *

3 75. Despite the swath of sociopolitical, technical, and logistical hurdles faced by social
4 media services, Meta's content moderation policies and enforcement systems are industry-leading,
5 comprehensive, and appropriately balanced as compared to the industry. To support this
6 conclusion, I offer an explicit benchmarking framework, as in the following. First, I discuss Meta's
7 iterative development of content moderation policies, which draw upon consultation with both
8 internal stakeholders and external experts. Second, I explain how these policies are operationalized
9 through Meta's moderation infrastructure, which combines machine learning classifiers, human
10 review, and user reports. I also compare Meta's performance to the leading industry standards of
11 the Digital Trust & Safety Partnership. In addition, I draw from my own academic and professional
12 experience conducting audits and evaluations of large-scale moderation systems to further
13 contextualize Meta's relative strengths.

14 **IV. Meta's Content Moderation Policies Are Comprehensive and Appropriately
15 Balanced**

16 **A. Overview and Development of Content Moderation Policy**

17 76. Meta has long been at the forefront of advances in content moderation policy and
18 systems to address and overcome, to the extent possible, the above-described challenges. As of
19 July 2023, Meta has spent roughly \$20 billion on user safety efforts since 2016, including on
20 content moderation.⁹⁰ Before describing Meta's content moderation systems, I first provide an
21 overview of Meta's comprehensive content moderation policies that direct the implementation of
22 its systems: the public-facing Community Standards, the internal Policy Labs, and internal
23 Borderline Content Policies, all of which govern Meta's content moderation systems.

24 77. Meta's policies are developed in conjunction with a panoply of internal and external

25
26 ⁸⁹ META3047MDL-001-00000112.

27 ⁹⁰ META3047MDL-113-00038326.

Highly Confidential (Competitor)

1 experts from around the globe who collaborate to develop objective policy rationales, to craft clear
2 descriptions of Meta's policies for notice to the public, and to help continually refine policies and
3 policy rationales over time as Meta enforces through its technologies.⁹¹ Meta also collaborates on
4 the operational side of content moderation, including by leveraging third-party signals and hashes,
5 such that different proprietary content moderation systems communicate information about
6 problematic content with each other and improve detection.⁹² This feedback loop helps to ensure
7 that Meta's policies continue to evolve with the realities of the technology on the ground.

8 78. Because of the amount of content on Meta's services, and because of Meta's careful
9 and thoughtful approaches to content moderation, total removal of all policy-violating content is
10 not possible. To get a sense of the scale, it is worth referring to Meta's latest available Full Year
11 Disclosure 2024,⁹³ which reports that the Family daily active people (DAP)—measure that Meta
12 utilizes to track how many active users are seen on Meta's family of apps on a daily basis—was
13 3.35 billion on average for December 2024. This was an increase of 5% year-over-year, suggesting
14 that the service's user base continues to grow. The size of Meta's user base, at almost four billion
15 active users, would make it the largest country in the world.⁹⁴ With this scale of a user base, there
16 are several, sometimes competing priorities that inform Meta's policy decisions: chiefly safety,
17 but also freedom of expression, privacy, and authenticity, among others.

18 79. Meta is uniquely situated to provide a service to foster free discussion,⁹⁵ and it

19
20 ⁹¹ See discussion *infra*.

21 ⁹² E.g., Meta. (2019, August). *Open-source photo and video matching*. Retrieved June 10, 2025,
22 from <https://about.fb.com/news/2019/08/open-source-photo-video-matching/>.

23 ⁹³ Meta. (2025). *Meta reports fourth quarter and full year 2024 results*. Retrieved June 10, 2025,
24 from <https://investor.atmeta.com/investor-news/press-release-details/2025/Meta-Reports-Fourth-Quarter-and-Full-Year-2024-Results/default.aspx>.

25 ⁹⁴ Meta. (2023). *Meta reports fourth quarter and full year 2023 results; Initiates quarterly dividend*. Retrieved June 10, 2025, from
26 https://s21.q4cdn.com/399680738/files/doc_financials/2023/q4/Meta-12-31-2023-Exhibit-99-1-FINAL.pdf.

27 ⁹⁵ Scholars have referred to Facebook as the place where people gather to discuss issues that are
28 vital to democracy. See, e.g. Heldt, A. P. M. (n.d.). *Merging the social and the public: How*

1 *Highly Confidential (Competitor)*

1 works to balance free discussion against its priority to provide “a safe environment for all who
2 engage in acts of expression.”⁹⁶ To encourage safe and free expression, Meta must follow a
3 principled approach when removing content.

4 80. To this end, Meta has developed and employed a baseline set of public-facing rules
5 known as the Community Standards that outline its general content moderation policies.⁹⁷ When
6 Meta updates these public-facing policies, it publishes a change log showing the updates made in
7 redline.⁹⁸

8 81. The Community Standards do not go into extreme detail as to what constitutes a
9 violation, and Meta does not otherwise disclose its internal guidelines that outline criteria for
10 content removal, downranking (i.e., demoting such content in news feeds), and other content
11 actions⁹⁹ because outlining its policy rationales and analytical frameworks for reviewing content
12 could provide bad actors with a guideline to make content that skirts its policies.¹⁰⁰ Instead, more
13 detailed internal guidelines, known as Policy Labs, lay out exactly what is and is not allowed on
14 Meta’s services.¹⁰¹ Meta develops and maintains a highly detailed Policy Lab for each category of
15 policy violating content relevant to this litigation. Each Policy Lab articulates a policy rationale,
16 identifies known or common questions that arise during human review, and sets operational

17 social media platforms could be a new public forum. In *Beyond the public square: Imagining digital democracy*. Even the Supreme Court has referred to social media as the “modern public square.” *Packingham v. North Carolina*, 137 S. Ct. 1730, 1737 (2017).

18 ⁹⁶ Meta. (n.d.). *Community Standards*. Retrieved June 10, 2025,
19 from <https://transparency.meta.com/policies/community-standards/>.

20 ⁹⁷ Meta. (n.d.). *Community Standards*. Retrieved June 10, 2025,
21 from <https://transparency.meta.com/policies/community-standards/>.

22 ⁹⁸ See, e.g., Meta. (n.d.). *Suicide and self-injury*. Retrieved June 10, 2025,
23 from <https://transparency.meta.com/policies/community-standards/suicide-self-injury/>.

24 ⁹⁹ See *infra* Part IV.

25 ¹⁰⁰ See, e.g., Meta. (2018, June). *Removing bad actors from Facebook*. Retrieved June 10, 2025,
26 from <https://about.fb.com/news/2018/06/removing-bad-actors-from-facebook/>.

27 ¹⁰¹ META3047MDL-177-00000003; META3047MDL-177-00000047; META3047MDL-177-
28 00000119; META3047MDL-177-00000170; META3047MDL-203-00203440 .

Highly Confidential (Competitor)

1 guidelines.¹⁰² Each Policy Lab also provides detailed examples of how certain content was
2 actioned. Further, Meta has an internal set of guidelines governing the review and enforcement
3 against “borderline content.”¹⁰³ Borderline content is that which “does not explicitly violate”
4 Meta’s Community Standards but which “may come close.”¹⁰⁴ Meta recognizes that universal
5 content policies come with inherent challenges, including that “some people have a definition of
6 what feels problematic that goes beyond what” Meta removes through its Community
7 Standards.¹⁰⁵ Getting the balance right requires “a range of expertise,” including policy, product,
8 and operations teams.¹⁰⁶

9 82. Internally, many Meta teams are focused on developing and refining content
10 moderation policy. Meta’s Content Policy, Integrity, Global Operations, and Research teams
11 collaborate to maintain Meta’s Policies to reflect the most up-to-date information with respect to
12 the parameters of harmful content and the relative risk it poses at any given time.¹⁰⁷ The Content
13 Policy team, which sits in more than a dozen locations around the world, is responsible for
14 developing Meta’s Community Standards and Community Guidelines.¹⁰⁸ The team includes
15 subject-matter experts across many harm types.¹⁰⁹ The Integrity team then assesses the global
16 impact of the potential policy change and builds the technology to scale the detection and
17

18 _____
19 ¹⁰² META3047MDL-177-00000003; META3047MDL-177-00000047; META3047MDL-177-
00000119; META3047MDL-177-00000170; META3047MDL-203-00203440.

20 ¹⁰³ META3047MDL-113-00000439.

21 ¹⁰⁴ META3047MDL-113-00000439.

22 ¹⁰⁵ META3047MDL-113-00000439.

23 ¹⁰⁶ META3047MDL-113-00000439.

24 ¹⁰⁷ Meta. (n.d.). *Deciding to change standards*. Retrieved June 10, 2025,
from <https://transparency.meta.com/policies/improving/deciding-to-change-standards/>.

25 ¹⁰⁸ Meta. (n.d.). *Deciding to change standards*. Retrieved June 10, 2025,
from <https://transparency.meta.com/policies/improving/deciding-to-change-standards/>.

26 ¹⁰⁹ Meta. (n.d.). *Deciding to change standards*. Retrieved June 10, 2025,
from <https://transparency.meta.com/policies/improving/deciding-to-change-standards/>.

Highly Confidential (Competitor)

enforcement of the updated policies.¹¹⁰ The Global Operations team (f/k/a/ the Community Operations team), which is responsible for partnering with outside partners and contractors to enforce Meta's Policies, keeps Meta informed of trends in enforcement so that Meta can either clarify its policies and/or upgrade its technology.¹¹¹ Further, Meta consults regularly with experts in adolescent development, psychology, and mental health to help ensure its services are safe and age-appropriate.¹¹²

At biweekly "Product Policy Forums," Meta convenes a variety of subject-matter experts from its "safety and cybersecurity policy teams, counterterrorism specialists, Global Operations employees, product managers, public policy leads and representatives from our legal, communications and diversity teams."¹¹³ Meta convenes such forums to discuss potential policy issues and propose policy changes, including but not limited to content moderation policies. In the interest of transparency, the minutes from each forum are published online.¹¹⁴ Policy proposals include perspectives gathered from discussions with academics, NGOs, or other experts and stakeholders outside the company. In the case of policies addressing SSI content, Meta has gathered perspectives from the National Eating Disorders Association, National Suicide Prevention Lifeline, Crisis Text Line, Samaritans, beyondblue, headspace, the Duke Center for Eating Disorders, victims, and others.¹¹⁵

¹¹⁰ Meta. (n.d.). *Deciding to change standards*. Retrieved June 10, 2025, from <https://transparency.meta.com/policies/improving/deciding-to-change-standards/>.

¹¹¹ Meta. (n.d.). *Deciding to change standards*. Retrieved June 10, 2025, from <https://transparency.meta.com/policies/improving/deciding-to-change-standards/>.

¹¹² Meta. (2024, January). *New protections to give teens more age-appropriate experiences on our apps*. Retrieved June 10, 2025, from <https://about.fb.com/news/2024/01/teen-protections-age-appropriate-experiences-on-our-apps/>.

¹¹³ Meta. (2018, November). *Product policy forum minutes*. Retrieved June 10, 2025, from <https://about.fb.com/news/2018/11/content-standards-forum-minutes/>.

¹¹⁴ Meta. (2018, November). *Product policy forum minutes*. Retrieved June 10, 2025, from <https://about.fb.com/news/2018/11/content-standards-forum-minutes/>.

¹¹⁵ META3047MDL-003-00081293; META3047MDL-003-00191358.

Highly Confidential (Competitor)

1 **B. Summary of Public-Facing Content Moderation Policies on Facebook or**
 2 **Instagram**

3 **1. Suicide & Self-Injury & Eating Disorder Content**

4 84. First, a summary of the key aspects of Meta's public-facing policies on suicide,
 5 self-injury, and eating disorder content is provided in the following policy card.¹¹⁶

6 *Table 1. Meta Policy Card on SSI & ED*

Policy Card	Details
Policy Name	Suicide, Self-Injury, and Eating Disorders
Policy Rationale	We care deeply about the safety of the people who use our apps. We regularly consult with experts in suicide, self-injury and eating disorders to help inform our policies and enforcement, and we work with organizations around the world to help people in distress. While we do not allow people to intentionally or unintentionally celebrate or promote suicide, self-injury or eating disorders, we do allow people to discuss these topics because we want our services to be a space where people can share their experiences, raise awareness about these issues, and seek support from one another.
Allowed Content	<ul style="list-style-type: none"> - Discussion of suicide, self-injury, or eating disorders aimed at sharing experiences, raising awareness, or seeking support. - Content about recovery that may contain non-graphic imagery (e.g., healed scars) placed behind a sensitivity screen.
Restricted Content	<ul style="list-style-type: none"> - Content that promotes, encourages, coordinates, or provides instructions for suicide, self-injury, or eating disorders. - Graphic self-injury imagery. - Depictions of a person who engaged in a suicide attempt or death by suicide. - Mocking victims or survivors of suicide, self-injury, or eating disorders. - Imagery depicting body modification in a suicide or self-injury context.
Content Labels and Restrictions	<ul style="list-style-type: none"> - Warning screens for sensitive content (e.g., healed cuts in a recovery context). - Age restriction (18 and older) for content depicting euthanasia/assisted suicide. - Regular consultation with external experts.
Enforcement	<ul style="list-style-type: none"> - Contacting emergency services for immediate risks.

24 85. In addition to its public-facing policies, Meta maintains an internal SSI & ED

25
 26 ¹¹⁶ Meta. (n.d.). *Suicide, self-injury, and eating disorders*. Meta Transparency Center. Retrieved
 27 May 4, 2025, from <https://transparency.meta.com/policies/community-standards/suicide-self-injury/>.

Highly Confidential (Competitor)

1 Policy Lab.¹¹⁷ While Meta has one Policy Lab that covers both SSI and ED content, these harm
2 types are addressed separately. In the SSI portions of this Policy Lab, Meta distinguishes between
3 different types of imagery depicting graphic suicide content. For example, Meta distinguishes
4 between imagery that depicts a suicide attempt or death by suicide, on the one hand, and imagery
5 that depicts a person engaging in euthanasia or physician-assisted suicide, on the other. While the
6 former would be removed, the latter would be allowed but placed behind a sensitivity screen
7 including a warning that the content may be disturbing. In the self-injury context, the Policy Lab
8 distinguishes between content that depicts graphic self-injury and content that depicts healed cuts
9 in the context of recovery.

10 86. To help ensure that its SSI & ED Policy Lab is robust and effective, Meta has
11 continuously consulted with mental health experts and organizations such as Forefront, Now
12 Matters Now, Save.org, and the National Suicide Prevention Lifeline.¹¹⁸ These consultations are
13 key because SSI content creates a unique conundrum in the context of content moderation: auto-
14 enforcement of all SSI content “could overenforce on people that are making cries for help” as
15 “[s]afety professionals mostly believe that SSI admission should provide resources to the users
16 and not punitive actions.¹¹⁹ Consistently, experts have told Meta that sometimes, certain SSI
17 content can be shared in a positive context and can help to destigmatize mental health struggles.¹²⁰
18 As explained by Dr. Reidenberg (Psy.D, FAPA, Executive Director, Save.org), whom Meta has
19 consulted in the development of its policy approach to SSI content:

20 _____
21 ¹¹⁷ META3047MDL-177-00000003.

22 ¹¹⁸ META3047MDL-001-00000112.

23 ¹¹⁹ META3047MDL-003-00079819.

24 ¹²⁰ META3047MDL-001-00000112 . One stakeholder explains, “Mental illness and thoughts of
25 suicide are just not something we talk about OPENLY. Yet talking and connecting is crucial to
26 helping prevent depression and suicide. The tools Facebook is rolling out, aim both at people
27 who are expressing suicidal thoughts and also guide concerned friends or family members to
resources and alternatives and appropriate interventions.” Meta. (n.d.). *Suicide prevention*.
Retrieved June 10, 2025,
from <https://about.meta.com/actions/safety/topics/wellbeing/suicideprevention>.

Highly Confidential (Competitor)

When someone expresses suicidal distress, it provides family, friends and even Facebook and Instagram the opportunity to intervene. If people can't share their pain, or it is shared and then removed, we've missed a chance to save someone's life. We train people to listen for this in conversations and to allow people to keep talking because we know that it is one way to help them through a crisis. Social media services allow us to do this in a way that brings many people to help very quickly.¹²¹

6 87. Furthermore, Meta recently collaborated with Snapchat on an initiative called
7 Thrive, through which the companies have shared SSI-related signals. Meta provided the technical
8 infrastructure for this initiative.¹²²

9 88. The SSI & ED Policy Lab establishes highly nuanced implementation standards for
10 ED content. These standards are designed to help ensure that violating content is actioned while at
11 the same time preserving the important role that Meta services play in promoting recovery and
12 connecting users with critical support. For example, Meta removes imagery depicting ribs,
13 collarbones, thigh gaps, hips, concave stomach, or protruding spine, unless that imagery is shared
14 in the context of recovery from an eating disorder.¹²³ But, when such content is shared in the
15 context of recovery, Meta places it behind a sensitivity screen rather than removing the content.¹²⁴

16 89. The SSI & ED Policy Lab also contains a glossary of terms that are relevant to
17 eating disorder content. For example, the glossary defines the term "long-shot" as "imagery where
18 part of the depicted person's shoulder or knee appears in the frame."¹²⁵ The Policy Lab also
19 addresses questions that Meta knows arise in the context of human review of eating disorder
20 content, such as which principles or criteria are considered when deciding whether a term or
21 concept is promoting or signaling eating disorders. The answers list hashtags that are associated

22 ¹²¹ Meta. (n.d.). *Suicide prevention*. Retrieved June 10, 2025,
23 from <https://about.meta.com/actions/safety/topics/wellbeing/suicideprevention>.

24 ¹²²Meta. (2024, September). *Preventing suicide and self-harm content from spreading online*.
25 <https://about.fb.com/news/2024/09/preventing-suicide-and-self-harm-content-spreading-online/>.

26 ¹²³ META3047MDL-177-00000003.

27 ¹²⁴ META3047MDL-177-00000003.

28 ¹²⁵ META3047MDL-177-00000003.

Highly Confidential (Competitor)

1 with eating disorders, such as “thinspo, bonespo, proana, promia, and thighgap.”¹²⁶

2. Bullying and Harassment Content

3 90. Next, a summary of the key aspects of Meta’s public-facing policies on bullying
4 and harassment (“BH”) is provided in the following policy card:

5 *Table 2. Meta Policy Card on Bullying & Harassment*¹²⁷

6 Policy Card	7 Details
Policy Name	Bullying and Harassment
Policy Rationale	Bullying and harassment happen in many places and come in many different forms: from making threats and releasing personally identifiable information to sending threatening messages and making unwanted malicious contact. We do not tolerate this kind of behavior because it prevents people from feeling safe and respected on Facebook, Instagram, and Threads. We distinguish between public figures and private individuals because we want to allow discussion, which often includes critical commentary of people who are featured in the news or who have a large public audience. For public figures, we remove attacks that are severe as well as certain attacks where the public figure is directly tagged in the post or comment. For private individuals, our protection goes further: We remove content that’s meant to degrade or shame, including, for example, claims about someone’s sexual activity. We recognize that bullying and harassment can have more of an emotional impact on minors, which is why our policies provide heightened protection for anyone under the age 18, regardless of user status.
Allowed Content	- Content that condemns or draws attention to bullying and harassment. - Discussion of public figures with critical commentary, if it does not include severe attacks or direct tagging.
Restricted Content	- Unwanted contact that is repeated, sexually harassing, or directed at many individuals with no prior solicitation. - Calls for self-injury or suicide. - Attacks based on experiences of sexual assault, exploitation, harassment, or domestic abuse. - Severe sexualized commentary and derogatory terms related to sexual activity. - Threats to release private information. - Content that degrades or expresses disgust towards individuals depicted in processes like menstruating, urinating, vomiting, or defecating.

24
25
26¹²⁶ META3047MDL-177-00000003.

27¹²⁷ Meta. (n.d.). *Bullying and harassment*. Meta Transparency Center. Retrieved May 4, 2025,
from <https://transparency.meta.com/policies/community-standards/bullying-harassment/>.

Highly Confidential (Competitor)

Policy Card	Details
Content	- Enhanced protections for minors under age 18.
Labels and Restrictions	- Warning screens for sensitive content.
Enforcement	- Regular consultation with external experts. - Bullying Prevention Hub for resources and support.

5 91. The BH Policy Lab establishes a series of “universal protections” that apply to all
 6 users, as well as more specific protections that depend on the type of user involved.¹²⁸ For example,
 7 Meta applies stricter protections when the user is a minor.¹²⁹ Similarly, Meta’s protections go
 8 further when the user is a private individual, as opposed to a public figure.¹³⁰

9 92. The BH Policy Lab is designed to address the context and intent behind a piece of
 10 content before actioning it. This helps Meta target bullying and harassment without inadvertently
 11 deleting content that was shared to condemn or draw attention to bullying and harassment, for
 12 example.

13 93. The first tier of protections, the universal protections, apply across all user types
 14 and across all types of content (e.g., organic content, paid/ads).¹³¹ These universal protections
 15 target, among other things, unwanted contact that is repeated, sexually harassing, or that calls for
 16 self-injury or suicide of a specific person or group. These universal protections also apply to
 17 content that threatens to release an individual’s private information.¹³²

18 94. The second tier of protections apply to all minors, private adults, and any limited-
 19 scope public figures.¹³³ Meta removes content that includes claims about sexual activity, except
 20 for claims that are in the context of criminal allegations against an adult. Meta also removes any
 21 content that highlights or otherwise draws negative attention to a person’s specific characteristics.

22 _____
 23 ¹²⁸ META3047MDL-177-00000047.

24 ¹²⁹ META3047MDL-177-00000047.

25 ¹³⁰ META3047MDL-177-00000047.

26 ¹³¹ META3047MDL-177-00000047.

27 ¹³² META3047MDL-177-00000047.

28 ¹³³ META3047MDL-177-00000047.

Highly Confidential (Competitor)

1 95. The third tier sets out additional protections for private minors, private adults, and
2 others.¹³⁴ Where this tier applies, Meta removes content such as claims about someone's sexual
3 orientation or content that expresses contempt or disgust about an individual. Many of these terms,
4 such as contempt or disgust, are defined in the Policy.

5 96. The fourth tier sets out additional protections for private minors, such as removing
6 any videos of physical bullying against a minor shared in a condemning context.¹³⁵

7 97. In addition to the tiers described above, the BH Policy Lab sets out specific
8 exceptions and allowances.¹³⁶ For example, though Meta would ordinarily remove videos of
9 physical bullying, it allows such videos in the context of fighting sports such as martial arts.
10 Furthermore, the Policy Lab sets out standards according to which Meta can identify the individual
11 being targeted in a bullying and harassing post. For example, identifying relevant responses from
12 the potential target in the comments when the parent post does not explicitly call out the target by
13 name or image; or the use of pronouns or the mentioning of an individual without use of their
14 name, along with comments indicating that the target referred to by the pronoun is known to the
15 people engaging with the post.¹³⁷ This is a complex task since bullies often do not name their target
16 but rather describe them (e.g., "this girl I work with . . .").

17 98. Finally, though the Policy Lab addresses a wide range of content, Meta also allows
18 users to submit reports to the service. This is particularly useful in the BH context because the
19 target may subjectively feel that they are being bullied or harassed, even if the content does not fit
20 into any of the categories outlined in the Policy.

21 **3. Adult Nudity & Sexual Activity Content**

22 99. A summary of the key aspects of Meta's public-facing policies on adult nudity and
23 sexual activity content is provided in the following policy card.

24

¹³⁴ META3047MDL-177-00000047.

25 ¹³⁵ META3047MDL-177-00000047.

26 ¹³⁶ META3047MDL-177-00000047.

27 ¹³⁷ META3047MDL-177-00000047.

Highly Confidential (Competitor)1 *Table 3. Meta Policy Card on ANSA*¹³⁸

Policy Card	Details
Policy Name	Adult Nudity and Sexual Activity
Policy Rationale	The display of nudity or sexual activity is restricted due to sensitivities related to cultural backgrounds and age. Exceptions are made for content intended for protest, awareness, educational, or medical reasons.
Allowed Content	<ul style="list-style-type: none"> - Images of female breasts in protest, breastfeeding, and post-mastectomy scarring. - Artistic depictions, such as paintings and sculptures, that include nude figures. - Medical or health-related imagery involving visible genitalia, anuses, or female nipples.
Restricted Content	<ul style="list-style-type: none"> - Photorealistic imagery of adult nudity, including visible genitalia, anuses, and female nipples, except in specific contexts. - Explicit sexual activity, including intercourse, oral sex, and the use of sex toys. - Implicit sexual activity unless in medical, health, or recognized fictional contexts. - Fetish-related activities that could lead to harm, including bestiality and incest. - Digital imagery of adult sexual activity unless it is for medical awareness, scientific discourse, or sexual health discussions.
Content Labels and Restrictions	<ul style="list-style-type: none"> - Sensitive content labels are added to certain images to inform viewers. - Viewing of some content is restricted to users aged 18 and older.
Enforcement	<ul style="list-style-type: none"> - Global review teams and external stakeholders help enforce and inform the policy. - Data on prevalence, content actioned, proactive rate, appealed content, and restored content is tracked and reported.

18 100. The ANSA Policy Lab also establishes universal protections that nevertheless give
 19 way to exceptions for nuanced situations.¹³⁹

20 101. ANSA content is prohibited outright in “ads and commerce surfaces.”¹⁴⁰ This
 21 includes “nudity, depictions of people in explicit or suggestive positions, or activities that are
 22

24 ¹³⁸ Meta. (n.d.). Adult nudity & sexual activity. Meta Transparency Center. Retrieved May 4,
 25 2025, from <https://transparency.meta.com/policies/community-standards/adult-nudity-sexual-activity/>.

26 ¹³⁹ META3047MDL-177-00000170.

27 ¹⁴⁰ META3047MDL-177-00000170.

Highly Confidential (Competitor)

1 overly sexually suggestive or provocative.”¹⁴¹

2 102. For organic (i.e., non-ad or non-commercial content), Meta “default[s] to removing
3 sexual imagery to prevent the sharing of non-consensual or underage content.”¹⁴²

4 103. In other contexts, “where appropriate and [] intent is clear,” Meta may allow ANSA
5 content shared “as a form of protest, to raise awareness about a cause or for educational or medical
6 reasons.”¹⁴³ Nevertheless, most ANSA content that is not defaulted to removal is filtered for users
7 under the age of eighteen or marked as sensitive.¹⁴⁴

8 **4. Violent & Graphic Content**

9 104. A summary of the key aspects of Meta’s public-facing policies on violent and
10 graphic (“VG”) content is provided in the following policy card.

11 *Table 4. Meta Policy Card on VG¹⁴⁵*

12 Policy Card	13 Details
13 Policy Name	14 Violent and Graphic Content
14 Policy Rationale	15 To protect users from disturbing imagery, we remove content that is particularly violent 16 or graphic, such as videos depicting dismemberment, visible innards or charred bodies. 17 We also remove content that contains sadistic remarks towards imagery depicting the 18 suffering of humans and animals. In the context of discussions about important issues 19 such as human rights abuses, armed conflicts or acts of terrorism, we allow graphic 20 content (with some limitations) to help people to condemn and raise awareness about 21 these situations.
21 Allowed Content	22 - Graphic content related to human rights abuses, armed conflicts, or acts of terrorism 23 with some limitations. 24 - Content with warning labels for graphic or violent imagery.

23 ¹⁴¹ META3047MDL-177-00000170.

24 ¹⁴² META3047MDL-177-00000170.

25 ¹⁴³ META3047MDL-177-00000170.

26 ¹⁴⁴ META3047MDL-177-00000170.

27 ¹⁴⁵ Meta. (n.d.). Violent & graphic content. Meta Transparency Center. Retrieved May 4, 2025,
from <https://transparency.meta.com/policies/community-standards/violent-graphic-content/>.

Highly Confidential (Competitor)

Policy Card	Details
1 Restricted Content	- Videos of dismemberment, visible internal organs, charred bodies, and victims of cannibalism.
2	- Sadistic remarks towards disturbing imagery.
3	- Live streams of capital punishment.
4	- Graphic imagery of animal suffering.
5 Content Labels and Restrictions	- Warning labels for graphic or violent imagery.
6	- Age restriction (18 and older) for certain graphic content.
7 Enforcement	- Regular consultation with external experts.
	- Specific measures for removing violent death videos upon request by family members.

8 105. Meta's VG Policy Lab protects users from sensitive content while also accounting
 9 for situations in which such content may serve a purpose, such as advocating against human rights
 10 abuses.¹⁴⁶ This Policy Lab is designed with the understanding that different people have different
 11 sensitivities to violent and graphic imagery. As such, Meta frequently adds warning screens to
 12 content instead of removing it entirely.

13 106. Furthermore, the VG Policy Lab accounts for situations in which graphic content
 14 may be acceptable.¹⁴⁷ For example, while content showing visible innards is generally not
 15 permitted, Meta does allow visible innards in the birthing context, though such content would be
 16 marked as sensitive.¹⁴⁸ Similarly, even if a video shows violent death or a life-threatening event,
 17 Meta may stop short of removing it because of the context of the video.¹⁴⁹ For example, Meta will
 18 mark as sensitive a video showing violence committed by uniformed police. And, while content
 19 depicting punctured skin may be removed in some contexts, it is generally allowed in the context
 20 of vaccines or piercings.¹⁵⁰

21 107. Meta has a slightly different set of policies for fictional graphic imagery, and these

23 ¹⁴⁶ META3047MDL-177-00000119.

24 ¹⁴⁷ META3047MDL-177-00000119.

25 ¹⁴⁸ META3047MDL-177-00000119.

26 ¹⁴⁹ META3047MDL-177-00000119.

27 ¹⁵⁰ META3047MDL-177-00000119.

Highly Confidential (Competitor)

policies distinguish between photorealistic and non-photorealistic imagery.¹⁵¹ Generally, Meta is more permissive of violent or graphic content when it is non-photorealistic.

5. Child Safety Content

108. A summary of the key aspects of Meta's public-facing policies on child safety content is provided in the following policy card.

Table 5. Meta Policy Card on

*Child Sexual Exploitation, Abuse, and Nudity (“Child Safety” or “CS”)*¹⁵²

Policy Name	Child Safety Content
Policy Rationale	We do not allow content or activity that sexually exploits or endangers children. When we become aware of apparent child exploitation, we report it to the National Center for Missing and Exploited Children (NCMEC), in compliance with applicable law. We know that sometimes people share nude images of their own children with good intentions; however, we generally remove these images because of the potential for abuse by others and to help avoid the possibility of other people reusing or misappropriating the images.
Allowed Content	Content from real-world art, cartoons, movies, or video games that depicts real or non-real non-sexual child abuse.
Restricted Content	<ul style="list-style-type: none">- Content, activity, or interactions that threaten, depict, praise, support, provide instructions for, make statements of intent, admit participation in, or share links of the sexual exploitation of children (including real minors, toddlers, or babies, or non-real depictions with a human likeness, such as in art, AI-generated content, fictional characters, dolls, etc.).- Content that solicits sexual content or activity depicting or involving children- Content that constitutes or facilitates inappropriate interactions with children- Content that attempts to exploit real children- Content (including photos, videos, real-world art, digital content, and verbal depictions) that sexualizes real or non-real children- Groups, Pages, and profiles dedicated to sexualizing real or non-real children

¹⁵¹ META3047MDL-177-00000119.

¹⁵² Child Sexual Exploitation, Abuse, and Nudity (n.d.). Meta Transparency Center. Retrieved on May 4, 2025 at <https://transparency.meta.com/policies/community-standards/child-sexual-exploitation-abuse-nudity/>.

Highly Confidential (Competitor)

Policy Name	Child Safety Content
	<ul style="list-style-type: none"> - Content that depicts real or non-real child nudity - Videos or photos that depict real or non-real non-sexual child abuse regardless of sharing intent, unless the imagery is from real-world art, cartoons, movies or video games - Content that praises, supports, promotes, advocates for, provides instructions for or encourages participation in non-sexual child abuse
Content Labels and Restrictions and Enforcement	<p>For the following content, we include a warning screen so that people are aware the content may be disturbing and limit the ability to view the content to adults ages eighteen and older:</p> <ul style="list-style-type: none"> - Videos or photos that depict police officers or military personnel committing non-sexual child abuse - Videos or photos of non-sexual child abuse, when law enforcement, child protection agencies, or trusted safety partners request that we leave the content on the service for the express purpose of bringing a child back to safety. <p>For the following content, we include a sensitivity screen so that people are aware the content may be upsetting to some:</p> <ul style="list-style-type: none"> - Videos or photos of violent immersion of a child in water in the context of religious rituals <p>For the following content, we include a warning label so that people are aware that the content may be sensitive:</p> <ul style="list-style-type: none"> - Imagery posted by a news agency that depicts child nudity in the context of famine, genocide, war crimes, or crimes against humanity, unless accompanied by a violating caption or shared in a violating context, in which case the content is removed

109. Reinforcing Meta's Community Standards, Meta's CS Policy Lab establishes universal protections that leave little room for exceptions and direct the removal of the overwhelming majority of CS content, with restrictions and limitations applied to the rest. The Policy Lab's Operational Guidelines, which human reviewers reference when reviewing potentially policy-violating content, contain highly granular factors for reviewers to consider.¹⁵³ For example, when determining minor age, reviewers are directed to consider explicit but also indirect information, such as the relative size of body parts, the relative roundness of faces, and the relative lack of a chin or jawline.¹⁵⁴

26 ¹⁵³ META3047MDL-203-00203440.

27 ¹⁵⁴ META3047MDL-203-00203440.

1 *Highly Confidential (Competitor)*

2 **C. Borderline Policies**

3 110. Further to its policies regarding clearly violating content, I understand Meta has
4 had internal “borderline content” policies since at least 2020.¹⁵⁵ Borderline Content Policies
5 address content that exists in a grey area, i.e., content that does not explicitly violate the
6 Community Standards but may still be considered harmful or controversial.¹⁵⁶ These policies are
7 designed to manage and mitigate the spread of such content, helping to ensure that the service
8 remains a safe and welcoming space for all users. Borderline content includes material that skirts
9 the edge of what is acceptable under Meta’s Community Standards.¹⁵⁷ This might include content
10 that is sensationalist, misleading, or provocative but not outright false or harmful.

11 111. For example, “Meta has an SSI Borderline Policy that covers things that are close
12 to prohibited, or SSI-adjacent like dark and depressing content.”¹⁵⁸ Meta’s ED Borderline Policy
13 is designed “to limit user exposure to content posted by non-connected users that could cause harm
14 when viewed either by a vulnerable user or as part of an aggregate viewing experience.”¹⁵⁹ Meta’s
15 VG Borderline Policy covers “things like pus or oozing wounds, people undergoing surgical
16 operations, fight videos, animal abuse, and fictional violence.”¹⁶⁰ Finally, Meta’s ANSA
17 Borderline Policy covers “things like implied sexual intercourse and stimulation.”¹⁶¹ Meta does
18 not have Borderline Policies for CS content because the Community Standards and CS Policy Lab

19

¹⁵⁵ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
20 30, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

21 ¹⁵⁶ META3047MDL-113-00000439.

22 ¹⁵⁷ META3047MDL-113-00000439.

23 ¹⁵⁸ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
24 31, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

25 ¹⁵⁹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
26 31, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

27 ¹⁶⁰ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
28 31, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

Highly Confidential (Competitor)

are the strictest set of content moderation policies and permit almost no, if any, nuance in the moderation of CS content.¹⁶²

112. Because borderline content is not policy violating, Meta does not remove borderline content but does apply various “soft” enforcement actions.¹⁶³ “Soft action” refers to an action that Meta takes on a piece of content short of removal to reduce its exposure to Meta’s users, such as downranking or demoting, filtering, age gating, adding warning screens or captions, non-recommendations, in-feed recommendation (IFR) filtering, or search engine results page (SERP) filtering.¹⁶⁴ Training machine learning models to detect borderline content is technically challenging and complex because it deals with content less amenable to categorization and assessment of the trade-offs between safety and voice.¹⁶⁵

V. Meta's Content Moderation Systems Are Robust and Reasonably Designed.

113. Meta utilizes a number of different systems to monitor and detect prohibited content and profiles on Facebook and Instagram. Enforcement is governed by three prongs:¹⁶⁶:

114. ***Removal:*** Meta removes content and, in certain instances, accounts that violate its policies. Users have an opportunity to appeal the removal.

115. ***Reduction:*** If content does not violate Meta’s content moderation policies but is still deemed sensitive, or borderline, Meta reduces its distribution and/or visibility generally or for a specific subset of users (e.g., users under the age of eighteen). Certain categories of content are barred (i.e., age gated) altogether from distribution or visibility to users under the age of eighteen. Meta also provides users with additional tools to opt in to certain automated content moderation features, such as hiding certain words.

¹⁶² See Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

¹⁶³ See *infra* Appendix A.

¹⁶⁴ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at 7, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

¹⁶⁵ See META 3047MDL-113-00000439.

¹⁶⁶ META 3047MDI -001-00000112

1 **Highly Confidential (Competitor)**

1 116. **Information:** Meta seeks to provide context on sensitive content before a user is
2 exposed to it (such as via warning screens).

3 117. This section proceeds in two parts. In the first, I will explain how Meta's robust
4 proactive detection and review mechanisms work and detail the specific mechanisms in place for
5 different types of harmful content (e.g., SSI, ED, BH, VG, ANSA, and CS). In the second part, I
6 will explain how Meta's network of human reviewers supports its proactive detection and review
7 mechanisms.

8 **A. Proactive Detection**

9 118. Meta identifies and actions most policy-violating content proactively, rather than
10 depending on users to report it. Meta uses AI tools to review and categorize content that may need
11 to be removed, downranked, or otherwise actioned.¹⁶⁷ One of the most important sets of tools are
12 "classifiers." Classifiers assess, detect, and categorize content, including text, video, and photos,
13 on both Facebook and Instagram to determine whether the content violates Meta's Community
14 Standards.¹⁶⁸ In addition, Meta has a suite of other automated features that help reduce user
15 exposure to prohibited content or abusive accounts.

16 **1. Meta's Early Leadership in Proactive Detection Technology**

17 **a. Facebook Immune System**

18 119. Meta has pioneered and iterated upon content moderation systems for many years,
19 beginning with the Facebook Immune System ("FIS"), launched by early 2011.¹⁶⁹ FIS was a first
20 solution for a real-time, large-scale integrity system and served as an early exemplar of adversarial
21 machine learning systems tailored for dynamic, high-volume environments. FIS was a cornerstone

22 ¹⁶⁷Meta. (n.d.). *How our enforcement technology works*. Retrieved June 10, 2025,
23 from <https://transparency.fb.com/enforcement/detecting-violations/how-enforcement-technology-works/>.

24 ¹⁶⁸ See *supra* Part IV; see also *infra* Appendix A (listing the topline classifiers used on Facebook
25 and Instagram as of February 2025).

26 ¹⁶⁹ Stein, T., Chen, E., & Mangla, K. (2011, April). *Facebook immune system*. In Proceedings of
27 the 4th workshop on social network systems.
<https://css.csail.mit.edu/6.858/2014/readings/facebook-immune.pdf>.

Highly Confidential (Competitor)

1 in Meta's defense infrastructure against adversarial threats such as phishing, spam, and abuse
2 across its social services.¹⁷⁰

3 120. At its core, FIS was a system designed to detect and mitigate threats such as spam,
4 fake accounts, misinformation, and other forms of abuse that can harm the integrity of the platform.
5 It did this by continuously monitoring and classifying billions of user interactions daily, including
6 reads and writes, to flag potential issues. It took automated actions such as removing content,
7 disabling accounts, or alerting human reviewers for further investigation to proactively defend the
8 social graph from abuse. FIS, initially designed to combat security threats, later evolved into a key
9 content moderation engine, enabling scalable, automated enforcement of policies related to
10 harmful content, misinformation, hate speech, and platform manipulation. FIS utilized a variety of
11 machine learning models and algorithms that analyze user behavior, content, and interactions to
12 identify suspicious activities. These models were continuously updated to adapt to new
13 threats. Further, FIS enabled real-time moderation of user-generated content with latencies under
14 50 milliseconds.¹⁷¹

15 121. Moderation decisions were informed by both explicit signals (e.g., "mark as spam")
16 and implicit behaviors (e.g., post deletions, rejected friend requests), which serve as training data
17 and anomaly triggers. Classifier services ran various ML algorithms (e.g., random forests, SVMs)
18 to evaluate content risk in real-time. Feature Extraction Language (FXL) enables on-the-fly
19 creation and deployment of new moderation features without system downtime. Policy Engine
20 translates classifier scores into enforcement actions (e.g., content takedown, account restriction)
21 via flexible, rule-based logic. Further, FIS aggregated features across multiple communication
22 channels (posts, messages, comments, friend requests) to detect violations that span modalities or

23
24 ¹⁷⁰ Stein, T., Chen, E., & Mangla, K. (2011, April). *Facebook immune system*. In Proceedings of
the 4th workshop on social network systems.

25 <https://css.csail.mit.edu/6.858/2014/readings/facebook-immune.pdf>.

26 ¹⁷¹ Stein, T., Chen, E., & Mangla, K. (2011, April). Facebook immune system (pp. 2–3). In
Proceedings of the 4th workshop on social network systems.

27 <https://css.csail.mit.edu/6.858/2014/readings/facebook-immune.pdf>

Highly Confidential (Competitor)

contexts.

122. In summary, FIS represented one of the earliest large-scale, adversarially aware content moderation systems and laid the groundwork for Meta's subsequent integrity and safety operations.

b. Research-Driven Technology

123. From a technological standpoint, pushing the state of the art in service integrity technologies is paramount. Not all services are equivalently engaged in the development of integrity systems and technologies. Some, but not all, services conduct state of the art research and publish scholarly work in various academic areas. Meta has been at the forefront for well over a decade in devoting substantial resources, efforts, and hiring of world known talent to study and address service integrity. Over a decade ago, for example, Meta instituted the Central Applied Science team (“CAS,” and formerly known as the “Core Data Science” team).¹⁷²

124. Meta has employed established academics to lead such research teams. For example, Dr. Lada Adamic, who until 2013 was a professor of computer science at the University of Michigan, served initially as manager then director of the CAS team from 2013 to 2020. Other academics have joined CAS, including Dr. Pablo Barberá (formerly a professor of political science at USC), Dr. Winter Mason (formerly a professor at Stevens Institute of Technology), and others.

125. The transparent publication and open documentation of service integrity protocols, technologies, and capabilities is another key element to advance public discourse and expert opinions. In this regard, not all social media service providers have engaged to the same extent with the public. Meta has a history of publishing highly impactful work documenting the internal

¹⁷² See Meta. (2025). *Academic Papers Citations*, Central Applied Science. Retrieved June 10, 2025, from https://scontent-iad3-2.xx.fbcdn.net/v/t39.2365-6/475662651_598728696446916_8597606275327940151_n.pdf?_nc_cat=105&ccb=1-7&_nc_sid=e280be&_nc_ohc=NACZKDy1AAIQ7kNvwGbeWiG&_nc_oc=AdkKRBxzatks2X30fVU4C2t8ve3ZB-Jyiqv01DKJi-9HImDZfNJObcRkVwFc1Wv697Y&_nc_zt=14&_nc_ht=scontent-iad3-2.xx&_nc_gid=I6cirBqvMSs0J0F3n4pQw&oh=00_AfPQKLbjafpiQFC135f8DeGHa3amSobo1kY6gcWC3TF_cQ&oe=686174C4 (showing scholarship published by CAS as early as 2009).

Highly Confidential (Competitor)

1 policies and systems that they design and deploy in production systems.

2 126. For example, a Facebook team led by Alon Halevy and collaborators, recently
 3 published a paper documenting the pipeline for enforcing integrity based on the guidelines
 4 developed at Meta. The paper appeared in the Communication of the ACM, the flagship computing
 5 magazine with a circulation of over 100,000 monthly readers.¹⁷³

6 127. Meta has been also providing support to the research community by starting
 7 initiatives to share (in a privacy preserving and ethical way), datasets to enable studies on ads
 8 targeting, content sharing, and engagement on Facebook.¹⁷⁴

9 **2. How Classifiers Work**

10 128. Building on years of technological iteration and research, Meta's suite of classifiers
 11 today are trained to look for specific signals and make predictions as to whether or not certain
 12 signals are present in a piece of content (e.g., nudity is present).¹⁷⁵ When a classifier is created,
 13 Meta trains the machine learning model by introducing more content of a specific type to improve
 14 the model's detection accuracy.¹⁷⁶

15 129. Teams of reviewers make the final call on classifier predictions, and classifiers
 16 learn from those decisions.¹⁷⁷ Meta measures recall during the pre- and post-deployment phases
 17 of classifier development to assess the technology's performance,¹⁷⁸ and continually refines the

18
 19 ¹⁷³ See Halevy, Alon, et al. "Preserving integrity in online social networks." *Communications of*
 the ACM 65.2 (2022): 92-98.

20 ¹⁷⁴ Meta. (n.d.). *Researcher datasets*. Retrieved June 10, 2025,
 from <https://fort.fb.com/researcher-datasets>.

21 ¹⁷⁵ Meta. (2024, Nov. 12). *How enforcement technology works*. Meta Transparency Center.
 <https://transparency.meta.com/enforcement/detecting-violations/how-enforcement-technology-works/>.

22 ¹⁷⁶ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at
 5-7, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

23 ¹⁷⁷ Meta. (2024, Nov. 12). *How enforcement technology works*. Meta Transparency Center.
 <https://transparency.meta.com/enforcement/detecting-violations/how-enforcement-technology-works/>.

24 ¹⁷⁸ Meta Defendants' First Supplemental Responses and Objections to Plaintiffs' Fifth Set of
 Interrogatories at 2-3, No. 4:22-MD-03047-YGR (N.D. Cal. Mar. 17, 2025).

1 *Highly Confidential (Competitor)*

1 models.¹⁷⁹

2 130. A classifier assigns a “Confidence Level” (“CL”) to the content it reviews.¹⁸⁰ A CL
3 is a statistical value referring to the probability that a piece of content is policy-violating (e.g., p-
4 value = 0.95); the higher the CL value, the more likely it is that the content is policy-violating.
5 Conversely, the lower the CL value, the less likely it is that the content is policy-violating.

6 131. CLs are measured against “threshold values” to determine what type of action, if
7 any, to take on a piece of content.¹⁸¹ If the CL exceeds the applicable threshold value, generally
8 the content will be automatically removed. If the CL is below the removal threshold, the classifiers
9 may nevertheless “soft action” the content or route it for human review, as described above in Part
10 IV.C.¹⁸² Soft actions may vary depending on the category of content and type of classifier.¹⁸³ If a
11 piece of content is assigned a CL that exceeds a certain threshold value, then that content is deemed
12 to violate Meta’s content policy and is automatically removed.¹⁸⁴ If a piece of content is assigned
13 a CL that is lower than the removal threshold value, but higher than the threshold value that would
14 trigger a soft action, then that content is routed for soft action and possible human review for a
15 final determination.¹⁸⁵ If the confidence level is relatively low, content is not sent for human
16

17
18

¹⁷⁹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
19 7, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

20 ¹⁸⁰ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
21 7, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

22 ¹⁸¹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
23 7, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

24 ¹⁸² Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
25 7, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

26 ¹⁸³ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
27 7-8, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

28 ¹⁸⁴ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
29 8, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

1 **Highly Confidential (Competitor)**

1 review; it is downranked, age gated, and/or subject to some other form of soft action.¹⁸⁶

2 132. Threshold values may vary by the type of policy-violating content (e.g., SSI, ED)
3 and may also vary across Facebook and Instagram. The threshold values are re-assessed both each
4 time classifiers are updated and on an ongoing basis in order to maintain a consistent accuracy bar.

5 133. Where a piece of content violates multiple Community Standards, and any of the
6 corresponding classifiers meet the threshold for deletion, the content will be removed.¹⁸⁷ Further,
7 multiple classifiers can trigger soft actions at the same time, leading to even stronger demotion.¹⁸⁸
8 Further, I understand there are certain “multi-class” classifiers that can identify overlapping policy
9 violations.¹⁸⁹

10 134. The AI detection system runs proactively and automatically, with technology
11 working behind the scenes to remove prohibited content, typically before any user sees it.¹⁹⁰ Meta
12 also continuously trains its artificial intelligence system for more accurate and nuanced detection
13 of prohibited content.¹⁹¹ One way Meta does so is by training classifiers to look at all of the
14 components of substantially similar posts, separate irrelevant differences from substantive
15 differences (e.g., the presence of a Web browser header in one image and its absence from a
16 substantially similar image would constitute an irrelevant difference), and then make a prediction

17

18¹⁸⁶ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
19 8, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

20¹⁸⁷ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
21 8, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

22¹⁸⁸ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
23 8, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

24¹⁸⁹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
25 8, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

26¹⁹⁰ Meta. (n.d.). *Prioritizing content review*. Retrieved June 10, 2025,
27 from <https://transparency.fb.com/policies/improving/prioritizing-content-review/>.

¹⁹¹Meta. (n.d.). *How our enforcement technology works*. Retrieved June 10, 2025,
from <https://transparency.fb.com/enforcement/detecting-violations/how-enforcement-technology-works/>.

Highly Confidential (Competitor)

1 of the contents' violative nature.¹⁹² Another way is by using a system that guides AI to cull training
2 data directly from millions of current content on Meta's services.¹⁹³ While this system also works
3 by reactively reviewing user reports, out of all policy-violating content removed on Facebook and
4 Instagram, the AI system generally removes over 90% proactively rather than through user
5 reporting.¹⁹⁴

6 135. Generally, if the item is obviously actionable (i.e., meets the threshold for deletion
7 or other soft action), the system is self-executing. Machine learning of this nature learns and
8 improves over time, and while it will never be perfect, its improvement has been observed in the
9 Community Standards and Enforcement Reports.¹⁹⁵

10 136. In addition to the classifiers discussed above, Meta also uses a banking system
11 whereby content that is flagged for human review is enqueued in the "Single Review Tool," or
12 SRT, to assess individual content, as described below in Part V.G. Policy-violating content is
13 uploaded to the Media Match Service ("MMS"), which identifies new and historical copies of the
14 content across Facebook and Instagram (and Messenger), which is automatically reviewed by
15 Meta's proactive AI system, comparing "new" content to the MMS bank.¹⁹⁶ Virtually every piece
16 of content uploaded to Facebook and Instagram (and Messenger) is checked against the MMS
17 bank.¹⁹⁷ Groups of content identified as copies can be auto-enforced en masse, or they can be
18 batched out for human review.¹⁹⁸ This en masse matching, auto-enforcement, and review process
19 helps improve capacity, as Meta can avoid rote review of duplicate pieces of content, and can
20

21 ¹⁹²See Meta. (n.d.). *How our enforcement technology works*. Retrieved June 10, 2025,
22 from <https://transparency.fb.com/enforcement/detecting-violations/how-enforcement-technology-works/>.

23 ¹⁹³ Meta. (n.d.). *Training AI to detect hate speech in the real world*. Retrieved June 10, 2025,
24 from <https://ai.meta.com/blog/training-ai-to-detect-hate-speech-in-the-real-world/>.

25 ¹⁹⁴ Antigone Davis Dep. Vol. 1 at 81:1-5, Mar. 4, 2025.

26 ¹⁹⁵ META3047MDL-001-00000112.

27 ¹⁹⁶ META3047MDL-050-00003108.

¹⁹⁷ See META3047MDL-050-00003108.

¹⁹⁸ META3047MDL-050-00003108.

Highly Confidential (Competitor)

quickly address violating recycled content multiple times. Meta has used banking since at least 2011.¹⁹⁹

137. Further, Meta maintains a “blackhole” database containing URLs that cannot be posted on Meta’s services.²⁰⁰

a. Content Moderation of Live (Real-Time) Content

138. Moderation of Live content (i.e., content streamed or broadcast in real time) on Facebook and Instagram also involves the use of classifiers. However, there is a relative lack of data associated with Live content (in contrast to other content) because Live is not as heavily utilized by users.²⁰¹ Enforcement on Live presents an additional challenge because the content is not static; classifiers have to be rerun every 30 seconds of Live video to reevaluate whether the video is violating.²⁰² In contrast, pre-recorded video uploaded to Meta's services is analyzed by a classifier just once.²⁰³ Furthermore, because of the nature of Live content, which can involve severe, rare, and immediate problems, moderation must be done near-instantaneously.²⁰⁴ Moderation of Live is also more complex because it requires its own custom classifiers.²⁰⁵

139. Live content classifiers are newer relative to their static-content counterparts, and they frequently result in content being enqueued for human review, except for a few high-severity

¹⁹⁹ META3047MDL-208-00050921 (linking an “overview of bank names since 2011”).

²⁰⁰ See META3047MDL-003-00170869; META3047MDL-003-00138204.

²⁰¹ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

²⁰² Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

²⁰³ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

²⁰⁴ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025); see also *infra* Appendix A.

²⁰⁵ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

Highly Confidential (Competitor)

1 harms.²⁰⁶ Generally, enqueueing happens at around a 60-65% confidence level, with autodeletion
 2 when classifiers are at approximately above 85-90%.²⁰⁷

3 **b. Abusive Account Moderation**

4 140. When an account has repeatedly posted prohibited content, Meta makes the account
 5 harder to find on its services and limits the recommendation of that account's posts. Pursuant to
 6 Meta's "strike policy," if a user has a certain number of qualifying content takedowns—typically
 7 within a year-long period—their account is disabled.²⁰⁸ Meta's strike policy is flexible, however,
 8 in that it may require fewer strikes before a takedown is actioned depending on the violation type—
 9 indeed, such action may occur in as little as a single post of prohibited content depending on the
 10 circumstances.²⁰⁹ For example, posting "very high severity+" content²¹⁰ typically results in
 11 automatic account disablement.²¹¹ When certain accounts are removed, Meta works to ensure that
 12 the device ID that created that account is blocked from creating new accounts in the app,²¹²
 13 particularly those posting "very high severity+" content.²¹³ Some things that are violations of
 14 Community Standards do not result in strikes (e.g., spam).²¹⁴

15 141. Meta automatically disables accounts that exhibit certain signals it monitors for

16 206 Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at
 17 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

18 207 Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at
 19 10-11, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

20 208 See META3047MDL-046-00005060; META3047MDL-003-00062744; META3047MDL-
 21 003-00038376.

22 209 See, e.g., META3047MDL-014-00240391.

23 210 META3047MDL-111-00372054. "Very high severity+" content is that which is "very rare,
 24 but can cause tremendous harm to" Meta's community and its reputation, and relevant categories
 25 of content include Child Safety and Credible Intent of Suicide (a subset of SSI) content. *Id.*

26 211 META3047MDL-014-00240391.

27 212 META3047MDL-004-00017504.

28 213 META3047MDL-208-00051129; see META3047MDL-014-00240391 (discussing blocking
 the device of a user allegedly engaged in child sextortion).

29 214 META3047MDL-046-00005060.

Highly Confidential (Competitor)

1 suspicious behavior. Meta identified and removed more than 90,000 accounts from August 1, 2023
2 to December 31, 2023 alone as a result of this monitoring.²¹⁵

3 142. Where accounts have bios that include certain prohibited terms, those accounts are
4 removed. For instance, an account with a bio that contains ED promotion (e.g., contains keywords
5 like “thinspo, bonespo, thinspiration, etc.”²¹⁶) will ordinarily be removed.²¹⁷

6 143. Abusive accounts are more resource intensive for Meta to review than individual
7 posts or comments because they can take almost twice the amount of time for humans to review.²¹⁸

8 **c. Content-Specific Classifiers**

9 144. Classifiers run automatically and review content on Facebook and Instagram. A
10 robust and reasonable content moderation system must have discernable categories of policy
11 violating content in order to scale effectively, and Meta’s content-specific classifiers, as described
12 below, offer such a solution. Different classifiers review and detect different categories of policy-
13 violating content. The following sections describe generally the classifiers used to detect,
14 categorize, and potentially action content across each of the relevant policy categories.

15 **d. Suicide and Self-Injury Classifiers**

16 145. I understand that Meta first implemented a classifier to action SSI content at least
17 as early as February 2019, and that Meta has continuously developed and refined multiple SSI
18 classifiers since then.²¹⁹ These classifiers use “pattern-recognition signals, such as phrases and
19 comments of concern, to identify possible distress.”²²⁰ But, detection of SSI content is incredibly
20 complex. For example, while humans might recognize that “I have so much homework I want to

21
22
23
24

²¹⁵ META3047MDL-034-00086282.

²¹⁶ META3047MDL-177-00000003.

²¹⁷ META3047MDL-003-00041618.

²¹⁸ META3047MDL-003-00070636.

²¹⁹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
11, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

²²⁰ Meta. (n.d.). *Suicide prevention*. Retrieved June 10, 2025, from
<https://about.meta.com/actions/safety/topics/wellbeing/suicideprevention>.

Highly Confidential (Competitor)

kill myself” is not a genuine cry for help, a computer does not have the contextual understanding of human nuance to not tag this phrase as alarming.²²¹ Similarly, a picture of train tracks may have to do with travel (and therefore be innocuous content), or it may relate to SSI.²²² Humans may have greater capacity than AI systems to more quickly tell the difference between these images based on the context, and training AI to do so is challenging, particularly at scale.²²³ To train a classifier to have this contextual understanding it needs to be fed both positive examples (what you want it to identify) and contrasting negative examples (what you do not want it to identify) so it can learn to distinguish the two.²²⁴ This requires a huge swath of data.²²⁵

146. Meta uses its AI technology, machine learning, and image-based technology (including pattern-recognition signals, such as phrases and comments of concern) to proactively identify and take action on clearly violating SSI content. SSI classifiers are refined by providing positive (i.e., actual SSI content) and negative (i.e., not actual SSI content, such as “I have so much homework, I want to kill myself”) examples, so it can learn to distinguish the two.²²⁶

147. Meta’s SSI classifiers review not only the content itself but also other information, for example, in the case of a post in which a user is seriously at risk, the classifier may detect and evaluate comments such as “Tell me where you are” or “has anyone heard from [user]?”²²⁷

148. Training automated technology to identify and distinguish between SSI content that

19 ²²¹ Facebook. (2018, September 10). *How Facebook AI helps suicide prevention.*
20 <https://about.fb.com/news/2018/09/inside-feed-suicide-prevention-and-ai/>.

22 ²²² META3047MDL-001-00000112.

23 ²²³ META3047MDL-001-00000112; *see supra* Part III.B.

25 ²²⁴ Meta. (2018, September). *How Facebook AI helps suicide prevention.*
26 <https://about.fb.com/news/2018/09/inside-feed-suicide-prevention-and-ai/>.

27 ²²⁵ Meta. (2018, September). *How Facebook AI helps suicide prevention.*
28 <https://about.fb.com/news/2018/09/inside-feed-suicide-prevention-and-ai/>.

1 **Highly Confidential (Competitor)**

1 is violative or non-violative is highly complex, particularly on Instagram, where much of the
2 content is image-based and therefore the intended meaning substantially relies on subtle context.²²⁸

3 **e. Eating Disorder Classifiers**

4 149. I understand Meta launched a classifier to action ED content specifically at least as
5 early as March 2021.²²⁹ Further, I understand Meta's ED classifiers have been continuously refined
6 since then. Prior to the launch of an ED-specific classifier, I understand ED content fell under the
7 purview of SSI policies and enforcement technology but were ultimately delineated into a unique
8 harm area because of the challenges unique to detecting and reviewing ED content. Because
9 Meta's services see less ED content than SSI content, there are fewer datapoints for the AI system
10 to learn from in the ED sphere than the SSI sphere. To combat this issue, Meta generates synthetic
11 training data the classifiers can learn from.

12 150. ED classifiers consider the following content as policy-violating and subject to
13 removal: “[c]ontent that focuses on depiction of ribs, collar bones, thigh gaps, hips, concave
14 stomach, or protruding spine or scapula when shared together with terms associated with eating
15 disorders.”²³⁰ Not all policy violating ED content is subject to removal, such as content that
16 “depicts ribs, collar bones, thigh gaps, hips, concave stomach, or protruding spine or scapula in a
17 recovery context.”²³¹ Rather, content of that nature is placed behind a sensitivity screen. Generally,
18 content promoting or encouraging EDs are removed, but users are allowed to post content that
19 touches on their own experiences around self-image and body acceptance.²³²

20 _____
21 ²²⁸ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at
22 12, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025)

23 ²²⁹ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at
24 15, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

25 ²³⁰ Meta. (n.d.). *Suicide and self-injury*. Retrieved June 10, 2025,
26 from <https://transparency.meta.com/policies/community-standards/suicide-self-injury/>.

27 ²³¹ Meta. (n.d.). *Suicide and self-injury*. Retrieved June 10, 2025,
from <https://transparency.meta.com/policies/community-standards/suicide-self-injury/>.

²³² Meta. (n.d.). *Suicide and self-injury*. Retrieved June 10, 2025,
from <https://transparency.meta.com/policies/community-standards/suicide-self-injury/>.

1 **Highly Confidential (Competitor)**

1 151. I understand that ED classifiers do not currently perform automatic deletions of
2 content; rather, content is either queued for human review and subsequent removal or else soft-
3 actioned.²³³

4 **3. Bullying and Harassment Classifiers**

5 152. I understand Meta first launched a classifier that actioned BH content at least as
6 early as May 2018.²³⁴ BH classifiers consider “severe” attacks on a public figure as policy-
7 violating, but not in all cases content that merely degrades or shames a public figure.²³⁵ Per Meta’s
8 policies, these classifiers consider public figures to be state and national level government officials,
9 political candidates for those offices, people with over one million fans or followers on social
10 media, and people who receive substantial news coverage. Further, BH classifiers consider content
11 meant to degrade or shame any private individual as policy-violating.²³⁶ Regardless of user status,
12 Meta recognizes that bullying and harassment can have more of an emotional impact on minors,
13 and extends heightened protection to anyone under the age 18.

14 153. Some content is removed for all users: repeatedly degrading or shaming contact,
15 attacks based on the target’s experience of sexual assault or domestic abuse, calls for SSI to a
16 specific individual or group of individuals, derogatory terms, claims that a violent tragedy did not
17 occur, claims that individuals are lying about being a victim of a violent tragedy, threats to release
18 private information, statements to engage in a sexual activity, severe sexualized commentary,
19 derogatory sexualized photoshop, calls for bullying or harassment of people, and degrading people
20 who are depicted vomiting, defecating, urinating, or menstruating.

21
22 _____
23 ²³³ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
24 16, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

24 ²³⁴ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
25 17, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

26 ²³⁵ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
27 18, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

28 ²³⁶ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
29 18, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

1 **Highly Confidential (Competitor)**

1 154. In addition to Meta’s tools to take action on a user’s behalf, Meta also provides
2 tools that help users control their own experiences on service, such as the ability to block certain
3 key words from being posted in the comments to their content.

4 **4. Graphic & Violent Content Classifiers**

5 155. I understand Meta first launched a classifier that actioned VG content at least as
6 early as December 2016.²³⁷ VG classifiers are intended to detect graphic and violent content on
7 Facebook and Instagram that may violate Meta’s Policies.²³⁸ VG classifiers consider, for example,
8 photos and videos of wounded or dead people if they show dismemberment, visible innards,
9 charred or burning people, victims of cannibalism, and/or throat slitting, as policy-violating.²³⁹

10 156. I understand that to protect users, and consistent with Meta’s Community
11 Standards, Meta employs a VG classifier to autodelete the most graphic and policy-violating
12 content.²⁴⁰ Other VG classifiers add a “disturbing” label to less graphic content to warn users that
13 the content may be sensitive or disturbing and allow users to disengage before viewing the
14 content.²⁴¹ When content is marked as disturbing by a VG classifier, the content is routed for
15 human review and potential deletion. Conversely, Meta allows a subset of less graphic content to
16 remain when users share content to shed light on or condemn acts.²⁴² Nevertheless, even for this
17 less graphic content, various VG classifiers “age gate” (i.e., age restrict) content from youth
18

19 ²³⁷ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
20 23, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

21 ²³⁸ Meta. (n.d.). *Violent & graphic content*. Retrieved May 4, 2025,
22 from <https://transparency.meta.com/policies/community-standards/violent-graphic-content/>.

23 ²³⁹ Meta. (n.d.). *Violent & graphic content*. Retrieved May 4, 2025,
24 from <https://transparency.meta.com/policies/community-standards/violent-graphic-content/>.

25 ²⁴⁰ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
26 23, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

27 ²⁴¹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
28 23, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

1 **Highly Confidential (Competitor)**

1 users.²⁴³

2 **5. Adult Nudity & Sexual Activity Classifiers**

3 157. I understand Meta first launched a classifier to detect and action ANSA content at
4 least as early as January 2016.²⁴⁴ This group of classifiers detect images and videos of adult nudity
5 and sexual activity on Facebook and Instagram that may violate policy.²⁴⁵

6 158. I understand adult nudity classifiers consider the following adult nudity content as
7 policy-violating: “visible genitalia”; “visible anuses and/or fully nude close-ups of buttocks”;
8 “uncovered female nipples”; “explicit sexual activity and stimulation”; “implicit sexual activity
9 and stimulation”; “erections”; “presence of by-products of sexual activity”; “sex toys placed upon
10 or inserted into mouth”; “stimulation of visible human nipples”; “squeezing female breasts”;
11 “imagery depicting fetish that involves [] acts that are likely to lead to the death of a person or
12 animal,” “dismemberment,” “cannibalism,” “feces, urine, spit, snot, menstruation or vomit,”
13 “bestiality,” and “incest”; “digital imagery of adult sexual activity”; and “extended audio of sexual
14 activity.”²⁴⁶

15 **6. Child Safety Classifiers**

16 159. Meta first launched a classifier to action “child safety” (“CS”) content at least as
17 early as 2018.²⁴⁷ CS classifiers proactively identify and action content and behaviors that may
18 violate Meta’s Community Standards on Child Sexual Exploitation, Abuse, and Nudity.²⁴⁸ (I

19 ²⁴³ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
20 23, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

21 ²⁴⁴ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
22 25, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

23 ²⁴⁵ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
24 25, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

25 ²⁴⁶ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
26 25, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

27 ²⁴⁷ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
28 27, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

Highly Confidential (Competitor)

1 understand that other experts are opining on Meta's child safety efforts more broadly, including
2 Meta's work to detect, remove, and report child sexual abuse material, or CSAM, which are outside
3 the scope of this report.)

4 160. Child Safety classifiers consider the following content as violating and subject to
5 removal, among others: content, activity, or interactions that threaten, depict, praise, support,
6 provide instructions for, make statements of intent, admit participation in, or share links of the
7 sexual exploitation of children; content that solicits sexual content or activity depicting or
8 involving children, nude imagery of real or non-real children, and/or sexualized imagery of real or
9 non-real children; content that solicits sexual encounters with children; content that constitutes or
10 facilitates inappropriate interactions with children; content that attempts to exploit real children by
11 coercing money, favors, or intimate imagery with threats to expose intimate imagery or
12 information or by sharing, threatening, or stating an intent to share private sexual conversations or
13 intimate imagery; content that sexualizes real or non-real children; Groups, Pages, and profiles
14 dedicated to sexualizing real or non-real children; content that depicts real or non-real child nudity;
15 videos or photos that depict real or non-real non-sexual child abuse regardless of sharing intent; or
16 content that praises, supports, promotes, advocates for, provides instructions for or encourages
17 participation in non-sexual child abuse.²⁴⁹

18 161. Further, I understand that when users report potential CS content or when CS
19 classifiers are triggered, the content is enqueued for human review, and Meta's review
20 prioritization system helps to ensure that the potential CS content most likely to violate policy
21 moves to the top of the review queue.²⁵⁰ I understand that this process is in addition to Meta's
22 other processes to detect, remove, and report CSAM, which are outside the scope of this report
23 and subject to the reports of other experts.

24
25

²⁴⁹ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at
26 27-28, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

27
²⁵⁰ Meta Defendants' Supplemental Responses and Objections to Plaintiffs' First Interrogatory at
28 28, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

Highly Confidential (Competitor)

1 **a. Borderline Classifiers**

2 162. Meta employs advanced AI and machine learning tools to detect and monitor
3 borderline content, which is, again, “content on Facebook and Instagram that is not strictly policy
4 implicating but nonetheless considered potentially problematic or low quality.”²⁵¹ These systems
5 analyze patterns and signals that indicate content may be approaching the threshold of
6 acceptability. Borderline classifiers complement Meta’s content-specific classifiers as described
7 above, balancing the technical benefit of discernable categories of policy violating content against
8 the ambitious goal of reducing exposure to content that is not clearly negative for all users.

9 163. A 2018 study published in *New Media & Society* showed how users actively discuss
10 strategies and successful approaches to circumventing Facebook content moderation guardrails,
11 often downplaying the significance of their intended actions.²⁵² Meta’s enforcement approach
12 generally involves reducing exposure to borderline content through soft actions rather than
13 removal, and some of those enforcement strategies are described below.

14 164. Due to the dynamic nature of these nuanced policies, continuous updates and
15 training of these AI models are necessary to keep up with evolving content trends and tactics used
16 to circumvent detection: approaches addressing challenges like the constantly changing violation
17 trends, lack of precision for specific violations, and limited exploration of new violation types.
18 Meta’s enforcement approach generally involves reducing exposure to borderline content through
19 soft actions rather than removal, and some of those enforcement strategies are described below.

20 **b. Content Downranking**

21 165. One of Meta’s primary strategies for handling borderline content is to reduce its
22 visibility. This can include “downranking,” i.e., demoting such content in news feeds, search

23
24 ²⁵¹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
25 30, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025). I understand that there are no
26 “borderline” content detection classifiers in the CS content category, *id.* at 28, as Meta intends its
27 policies to address CS content as directly as possible.

28 ²⁵² Gerrard, Y. (2018). *Beyond the hashtag: Circumventing content moderation on social media.*
New Media & Society, 20(12), 4492-4511.

Highly Confidential (Competitor)

1 results, and recommendations to limit its reach without outright removal. By making borderline
2 content less prominent, Meta aims to reduce its potential impact while avoiding direct censorship.

3 166. A recent academic study, published by researchers at the Swiss' EPFL in
4 collaboration with a Meta scientist, showed that automated content moderation (comment deletion)
5 decreased subsequent rule-breaking behavior.²⁵³ Furthermore, the study found that "the effect of
6 deletion on the affected user's subsequent rule-breaking behavior was longer-lived than its effect
7 on reducing commenting in general, suggesting that users were deterred from rule-breaking but
8 not from commenting." I find Meta's automated approach to reduce the visibility of certain
9 borderline content quite reasonable and robust, as supported by the literature, as a means to
10 decrease subsequent posting of such content without deterring users from expressing themselves.

11 **c. User Warning Screens and Labels**

12 167. Providing users with context through warning screens and labels, informing users
13 of the nature of the content that a user must acknowledge before they can view it, can also help to
14 mitigate the effects of borderline content, to alert users to potentially sensitive content, and to allow
15 users to disengage with the post should they choose. For example, posts related to suicide or suicide
16 attempts that are not policy-violating receive warning screens, as do violent and graphic posts,
17 some forms of nudity, and posts raising awareness on bullying and harassment.²⁵⁴ These labels
18 and warning screens empower users to make informed decisions about the content they engage
19 with and share.

20 **B. Additional Automated Tools to Prevent Exposure to Harmful Content**

21 168. The "Hidden Words" tool allows a user to set up his or her account such that posts

22
23 ²⁵³ Horta Ribeiro, M., Cheng, J., & West, R. (2023, April). Automated content moderation
24 increases adherence to community guidelines. In Proceedings of the ACM web conference 2023
(pp. 2666-2676).

25 ²⁵⁴ Meta. (2025, Mar. 13). *Providing context on sensitive or misleading content*. Meta
26 Transparency Center. <https://transparency.meta.com/enforcement/taking-action/context-on->
27 sensitive-misleading-content/ (discussing examples on Facebook); Meta. (n.d.). *Suicide and self-injury*. Retrieved June 10, 2025, from <https://transparency.meta.com/policies/community-standards/suicide-self-injury/> (discussing similar warning screens on Instagram).

Highly Confidential (Competitor)

that contain certain words can be hidden from his or her feed on Instagram. The same feature applies to direct messages as well; users can turn on this feature to automatically filter direct message requests containing offensive words, phrases, and emojis,²⁵⁵ placing them in a folder separate from innocuous messages so that users are shielded from negative messages unless they opt in.²⁵⁶ Once a user opens the folder with the filtered messages, the message text is still automatically blurred until the user taps to uncover it.²⁵⁷ The program was launched in April 2021.²⁵⁸ Instagram worked with outside experts to develop a preset list of terms that will be filtered from DMs, and the user can also create a custom list of keywords that will be filtered out.²⁵⁹

169. “On-Device Nudity Control” (“ODNC”) is a tool that detects unwanted nude content received in Instagram direct messages and automatically blurs it out. ODNC signals to the user that an image may contain nudity, and gives them the option to unblur the media, block the message, or get further help.²⁶⁰ It is controlled by the user, and does not automatically report images back to Meta or law enforcement authorities.²⁶¹ It is on by default for users under 18

²⁵⁵ Instagram. (n.d.). Hide comments or message requests you don't want to see on Instagram. Retrieved June 10, 2025, from <https://help.instagram.com/700284123459336>.

²⁵⁶ META3047MDL-003-00001792; META3047MDL-053-00012559.

²⁵⁷ Instagram. (n.d.). *Introducing new tools to protect our community from abuse*. Retrieved June 10, 2025, from <https://about.instagram.com/blog/announcements/introducing-new-tools-to-protect-our-community-from-abuse>.

²⁵⁸ Meta. (n.d.). *Our tools, features, and resources to help support teens and parents*. Retrieved June 10, 2025, from <https://www.meta.com/help/policies/safety/tools-support-teens-parents/>.

²⁵⁹ Instagram. (n.d.). *Introducing new tools to protect our community from abuse*. Retrieved June 10, 2025, from <https://about.instagram.com/blog/announcements/introducing-new-tools-to-protect-our-community-from-abuse>.

²⁶⁰ Meta. (2024, April). *New tools to help protect against sextortion and intimate image abuse*. <https://about.fb.com/news/2024/04/new-tools-to-help-protect-against-sextortion-and-intimate-image-abuse/>.

²⁶¹ Meta. (2024, April). *New tools to help protect against sextortion and intimate image abuse*. <https://about.fb.com/news/2024/04/new-tools-to-help-protect-against-sextortion-and-intimate-image-abuse/>.

1 **Highly Confidential (Competitor)**

1 globally, and opt-in for everyone else.²⁶²

2 170. “Comment Filters/Comment Cover” are tools which block toxic, offensive, and
3 divisive comments automatically.²⁶³ The feature was first launched on Instagram in June 2017. It
4 was expanded in May 2018 to include bullying and offensive comments, automatically hiding
5 toxic comments, comments attacking a person’s appearance or character.²⁶⁴ Upon its
6 implementation, Meta saw an 8% drop in bullying comment reports.²⁶⁵ However, this feature
7 experienced some shortcomings, such as difficulty recognizing foreign characters.²⁶⁶ The
8 comment filter is on by default, but users can turn it off in the Comment Controls center in the
9 Instagram app. Additionally, there is a manual comment filter that enables the user to set their own
10 words that trigger the filter.²⁶⁷ Meta uses “fuzzy matching” to catch variations of prohibited
11 words.²⁶⁸

12 171. “Sensitive Content Control” is a tool unique to Instagram which allows users to
13 decide how much sensitive content shows up in Instagram’s Explore tab, was launched in July
14

15 _____
16 ²⁶² Meta. (2024, April). *New tools to help protect against sextortion and intimate image abuse.*
17 from <https://about.fb.com/news/2024/04/new-tools-to-help-protect-against-sextortion-and-intimate-image-abuse/>.

18 ²⁶³ Instagram. (2017, June 29). *Keeping Instagram a safe place for self-expression.*
<https://about.instagram.com/blog/announcements/keeping-instagram-a-safe-place-for-self-expression>.

19 ²⁶⁴ Instagram. (n.d.). *Bully Filter and Kindness Prom to protect Instagram community.* Retrieved
20 June 10, 2025, from <https://about.instagram.com/blog/announcements/bully-filter-and-kindness-prom-to-protect-our-community>; Meta. (n.d.). *Our tools, features, and resources to help support teens and parents.* Retrieved June 10, 2025,
21 from <https://www.meta.com/help/policies/safety/tools-support-teens-parents/>.

22 ²⁶⁵ META3047MDL-003-00139760.

23 ²⁶⁶ META3047MDL-003-00058255.

24 ²⁶⁷ Instagram. (2017, June 29). *Keeping Instagram a safe place for self-expression.* Retrieved
25 from <https://about.instagram.com/blog/announcements/keeping-instagram-a-safe-place-for-self-expression>.

26 ²⁶⁸ For instance, if “taco” is prohibited, through fuzzy matching, taaco, ta-co, and other variants
27 would be caught as well. META3047MDL-003-00046768.

Highly Confidential (Competitor)

1 2021.²⁶⁹ Sensitive content covers posts that are not against Meta's policies but may be upsetting
 2 to certain users (e.g., sexually suggestive or violent content).²⁷⁰ Sensitive Content Control gives
 3 users control over the strictness of the filter with three options: More (allows users to see more
 4 sensitive content), Standard (default state), and Less (restricts more sensitive content).²⁷¹ Users
 5 under the age of 18 cannot opt into More sensitive content, as only the Standard and Less options
 6 are available.²⁷²

7 **1. Public Efforts to Improve Content Moderation Technology**

8 172. In recognition of the challenges of detecting harmful content, Meta open-sources
 9 some of its technology to make it available for others to use and improve.²⁷³ For other content
 10 moderators who use their own content-matching technology, hash-sharing enables those systems
 11 to share digital fingerprints of problematic content to make detection easier.²⁷⁴ Meta has also
 12 launched open competitions, in partnership with Microsoft, the Partnership on AI, academics from
 13 several universities, Getty Images, and Driven data, to attract talent to solve difficulties around

14 269 Instagram. (n.d.). *Introducing new tools to protect our community from abuse*. Retrieved June
 15 10, 2025, from <https://about.instagram.com/blog/announcements/introducing-new-tools-to-protect-our-community-from-abuse>. As of a year after implementing, 99% of users who were
 16 defaulted in remain with this setting on. Meta. (2024, January). *Our work to help provide young
 17 people with safe, positive experiences*. Retrieved June 10, 2025,
 18 from <https://about.fb.com/news/2024/01/our-work-to-help-provide-young-people-with-safe-positive-experiences/>.

19 270 Instagram. (n.d.). *Introducing new tools to protect our community from abuse*. Retrieved June
 20 10, 2025, from <https://about.instagram.com/blog/announcements/introducing-new-tools-to-protect-our-community-from-abuse>.

21 271 Instagram. (n.d.). *Introducing new tools to protect our community from abuse*. Retrieved June
 22 10, 2025, from <https://about.instagram.com/blog/announcements/introducing-new-tools-to-protect-our-community-from-abuse>.

23 272 META3047MDL-014-00243437.

24 273 See, e.g., Meta. (2019, Sept. 5). *Creating a dataset and a challenge for deepfakes*. Meta Blog.
 25 <https://ai.facebook.com/blog/deepfake-detection-challenge/>; Meta. (2019, August). *Open-
 26 sourcing photo- and video-matching technology to make the internet
 27 safer*. <https://about.fb.com/news/2019/08/open-source-photo-video-matching/>.

28 274 Facebook. (2019, August). *Open-sourcing photo- and video-matching technology to make the
 29 internet safer*. <https://about.fb.com/news/2019/08/open-source-photo-video-matching/>.

Highly Confidential (Competitor)

1 detecting deepfakes and hateful memes. Two competitions that I am aware of were called the
2 “Deepfake Detection Challenge,”²⁷⁵ and the “Hateful Memes Challenge.”²⁷⁶

3 173. *Deepfake Detection Challenge*: The Deepfake Detection Challenge was intended
4 to “produce technology that everyone can use to better detect when AI has been used to alter a
5 video in order to mislead the viewer.”²⁷⁷ It includes monetary grants and awards to motivate the
6 industry to innovate in this space, and is overseen by “the Partnership on AI’s new Steering
7 Committee on AI and Media Integrity,” which is made up of a broad cross-sector coalition of
8 organizations including Meta (then-Facebook), WITNESS, Microsoft, and others in civil society
9 and the technology, media, and academic communities.²⁷⁸ Meta invested heavily in this challenge
10 (over \$10 million), through which it has tried to create as realistic a dataset as possible in order to
11 get the most useful results.²⁷⁹

12 174. *Hateful Meme Challenge*: Similarly, the Hateful Memes challenge was a first-of-
13 its-kind \$100,000 competition intended to accelerate research on the problem of detecting hate
14 speech that combines both images and text.²⁸⁰ It had more than 3,300 participants around the
15 world, and the winning model achieved a 0.8450 AUC ROC, far exceeding the baseline model.²⁸¹

16
17
18 275 Meta. (2019, Sept. 5). *Creating a dataset and a challenge for deepfakes*. Meta Blog.
19 <https://ai.facebook.com/blog/deepfake-detection-challenge/>.

20 276 Meta. (2020, Dec. 11). *Hateful memes challenge winners*. Meta Blog.
21 <https://ai.facebook.com/blog/hateful-memes-challenge-winners/>.

22 277 Meta. (2019, Sept. 5). *Creating a dataset and a challenge for deepfakes*. Meta Blog.
23 <https://ai.facebook.com/blog/deepfake-detection-challenge/>.

24 278 Meta. (2019, Sept. 5). *Creating a dataset and a challenge for deepfakes*. Meta Blog.
25 <https://ai.facebook.com/blog/deepfake-detection-challenge/>.

26 279 Meta. (2019, Sept. 5). *Creating a dataset and a challenge for deepfakes*. Meta Blog.
27 <https://ai.facebook.com/blog/deepfake-detection-challenge/>.

28 280 Meta. (2020, Dec. 11). *Hateful memes challenge winners*. Meta Blog.
29 <https://ai.facebook.com/blog/hateful-memes-challenge-winners/>.

Highly Confidential (Competitor)

1 **2. Other Efforts to Improve Content Moderation Technology**

2 175. Models developed by Meta AI scientists, such as XLM-R,²⁸² leverage cross-lingual
 3 training to understand and moderate content across different languages, enabling them to apply
 4 knowledge gained from one language to another.

5 176. Meta develops and employs advanced AI transformer models, such as the
 6 RoBERTa²⁸³ and XLM-R²⁸⁴ architectures. Transformer models such as these are at the forefront
 7 of detecting complex and subtle forms of harmful content. These models rely on mechanisms to
 8 focus on relevant parts of the input data, allowing them to understand the context and nuances of
 9 language and imagery.

10 177. Another notable advancement is the Linformer architecture, developed by Meta's
 11 scientists, which enhances the efficiency of Transformer models, enabling the processing of longer
 12 text inputs with less computational overhead.²⁸⁵

13 **C. Meta's Reporting System**

14 178. While Meta's proactive mechanisms identify the bulk of the harmful content posted
 15 to Meta's services, given the immense volume and highly varied nature of the content on Meta's
 16 services, no technology could ensure that 100% of all policy-violating content is appropriately
 17 actioned. Thus, Meta has a "safety net" of reactive content moderation that can review and take

21 ²⁸² See Conneau, Alexis, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume
 19 Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin
 20 Stoyanov. *Unsupervised cross-lingual representation learning at scale*. arXiv preprint
 21 arXiv:1911.02116 (2019). Note that this paper has been cited over 7,000 times by researchers in
 the field of AI, demonstrating its significant impact.

22 ²⁸³ See Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., ... & Stoyanov, V. (2019).
 Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.

23 ²⁸⁴ See Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., ... &
 24 Stoyanov, V. (2019). *Unsupervised cross-lingual representation learning at scale*. arXiv preprint
 arXiv:1911.02116

25 ²⁸⁵ See Wang, Sinong, Belinda Z. Li, Madian Khabsa, Han Fang, and Hao Ma. *Linformer: Self-
 26 attention with linear complexity*. arXiv preprint arXiv:2006.04768 (2020). Note that this paper
 27 has been cited over a thousand times by researchers in the field of AI, substantiating its
 contribution and importance.

Highly Confidential (Competitor)

1 action on content that automated systems did not. This includes content that is reported by users.
 2 Instagram alone receives roughly 10 million reports every day.²⁸⁶

3 179. Adding to the challenge of scale, user reports often lack reliable indicators of the
 4 actual saliency of policy violating content. The user reporting feature has been known to be abused:
 5 Academic research has shown that reporting tools are often used adversarially rather than to flag
 6 legitimate violations.²⁸⁷ For example, “mass reporting brigades” are common on Twitter/X²⁸⁸ and
 7 Reddit,²⁸⁹ as well as Meta’s services.²⁹⁰ Further, ideological weaponization of reporting systems
 8 disproportionately impacts marginalized users.²⁹¹ Services can struggle to identify malicious use
 9 of moderation tools, especially when decisions are automated or unverified.

10 180. In addition, a Cornell researcher recently postulated the notion that users have been
 11 known to make use of higher-severity issue tagging for lower-severity violations in their reports
 12 in hopes of expediting response times or engage in harassment.²⁹² Reports are therefore not always
 13 an accurate signal of a potential violation type, which Meta’s systems take into account.

14 _____
 15 ²⁸⁶ META3047MDL-019-00093833.

16 15 ²⁸⁷ Myers West, S. (2018). *Censored, suspended, shadowbanned: User interpretations of content*
 17 *moderation on social media platforms*. New Media & Society, 20(11), 4366-4383; Gorwa, R.,
 18 Binns, R., & Katzenbach, C. (2020). *Algorithmic content moderation: Technical and political*
challenges in the automation of platform governance. Big Data & Society, 7(1),
 2053951719897945.

19 19 ²⁸⁸ Zannettou, S., Caulfield, T., Setzer, W., Sirivianos, M., Stringhini, G., & Blackburn, J. (2019,
 20 June). *Who let the trolls out? towards understanding state-sponsored trolls*. In Proceedings of
 the 10th ACM Conference on Web Science (pp. 353-362).

21 21 ²⁸⁹ Jhaver, S., Appling, D. S., Gilbert, E., & Bruckman, A. (2019). “Did you suspect the post
 22 would be removed?” Understanding user reactions to content removals on Reddit. Proceedings
 of the ACM on human-computer interaction, 3(CSCW), 1-33.

23 23 ²⁹⁰ See, e.g., Meta. (2021, December). Meta’s adversarial threat report. Retrieved June 10, 2025,
 24 from <https://about.fb.com/news/2021/12/metas-adversarial-threat-report/#:~:text=Mass%20Reporting%3A%20We%20will%20remove,Behavior%20policy%20against%20mass%20reporting>.

25 25 ²⁹¹ Marwick, A. E., & Caplan, R. (2018). *Drinking male tears: Language, the manosphere, and*
networked harassment. Feminist Media Studies, 18(4), 543-559.

26 26 ²⁹² Meisner, C. (2023). *The weaponization of platform governance: Mass reporting and*
algorithmic punishments in the creator economy. Policy & Internet, 15(4), 466-477.

Highly Confidential (Competitor)

1 181. As noted elsewhere in this report, users are typically not exposed to very much
2 prohibited content.²⁹³ In the rare instance where potentially prohibited content is not filtered by
3 the AI system, users are able to report such content through in-app reporting functions, which are
4 available on all devices.²⁹⁴

5 182. Meta employs metrics to assess and continually improve action taken on user
6 reports, such as Recall over All Reports (“RoAR”), which was replaced by Recall over Tagged
7 Reports (“RoTR”) in 2024.²⁹⁵ RoAR asked what percentage of user reports on actually violating
8 content did Meta enforce within N days.²⁹⁶ RoTR, the evolution of the RoAR metric, measures the
9 percentage of reported and confirmed “very high severity+” content that is appropriately actioned
10 within 96 hours.²⁹⁷

11 **D. Meta’s Human Review Processes**

12 183. Meta has thousands of reviewers globally who review potential violations of
13 Facebook and Instagram policies, which is called the Global Operations team.²⁹⁸ I will first provide
14 an overview of the review tool, the review process, and escalation discuss how Meta undertakes
15 this effort.

16 **1. Overview**

17 184. Meta’s human reviewers use the “Single Review Tool” or SRT to assess individual

18 293 See *infra* Part VI (discussing the prevalence rate in Meta’s CSER report); Antigone Davis
19 Dep. Vol. 1 at 81:1-5, Mar. 4, 2025.

20 294 Instagram. (n.d.). *Report a post or profile on Instagram*. Instagram Help Center. Retrieved
21 June 10, 2025, from <https://help.instagram.com/192435014247952/>; Meta. (2023, Oct. 18). *How*
22 *technology detects violations*. Meta Transparency Center.
<https://transparency.meta.com/enforcement/detecting-violations/technology-detects-violations/>;
META3047MDL-001-00000112.

23 295 META3047MDL-208-00061906.

24 296 META3047MDL-208-00050989.

25 297 META3047MDL-208-00061906; see *supra* footnote 210 (defining “very high severity+”
content).

26 298 Meta. (n.d.). *How review teams work*. Meta Transparency Center. Retrieved June 10, 2025,
27 from <https://transparency.meta.com/enforcement/detecting-violations/how-review-teams-work/>.

Highly Confidential (Competitor)

1 content. What reviewers see depends on what is reported. For instance, where Instagram comments
2 are reported, the reviewer sees the Instagram image, the caption, and the comment under review.²⁹⁹
3 They then can determine whether the comment complies with the Community Standards and
4 Policy Labs.³⁰⁰ At times, the reviewer simply looks at the content to determine if it is violates
5 Meta's policies, but at other times, a reviewer may be encouraged to look at the context of the
6 content before reaching a decision.³⁰¹

7 185. If a human reviewer needs help, they can escalate the questionable content to a
8 subject matter expert (i.e., someone specialized in the review of a particular category of content)
9 to decide.³⁰² Additionally, where a reviewer is concerned about an imminent credible threat of
10 physical harm, the content is escalated to safety experts who determine the threat's credibility and
11 whether it should be reported to emergency services.³⁰³

12 186. Meta conducts regular audits and quality checks to help ensure that reviewers are
13 applying standards consistently and accurately.³⁰⁴ These audits help identify areas for
14 improvement in training and guidelines.

15 187. Meta also has a system in place to allow for two kinds of appeals—"actor appeals,"
16 where a user appeals the decision to take down something they themselves posted; and "reporter
17 appeals" where the reporter asks Meta to "seriously, take a look" at something that was *not*
18 removed.³⁰⁵ Similar to first-level user reports, appeals are reviewed by Meta's automated systems
19 for review prioritization.

20
21 ²⁹⁹ META3047MDL-003-00036694.

22 ³⁰⁰ META3047MDL-003-00036694.

23 ³⁰¹ META3047MDL-003-00036694.

24 ³⁰² META3047MDL-001-00000112 .

25 ³⁰³ META3047MDL-001-00000112 .

26 ³⁰⁴ META3047MDL-003-00001478; META3047MDL-050-00084511.

27 ³⁰⁵ META3047MDL-003-00144080. By at least Feb. 11, 2019, actor appeals were enabled for
most problem types, and reporter appeals were enabled across the board.

1 **Highly Confidential (Competitor)**

2 **2. In-House & Third-Party Human Reviewers**

3 188. Meta has around 35,000 internal and third-party human reviewers. Internal
4 reviewers are closely-integrated with Meta's overall content moderation strategy. This integration
5 helps ensure that content moderation is consistent with company policies and that human reviewers
6 are aligned with the company's mission and standards.

7 189. Internal reviewers can provide immediate feedback on policy effectiveness and
8 suggest improvements. This direct line of communication helps in refining content policies and
9 enforcement mechanisms, leading to more accurate enforcement outcomes.

10 190. Inviting third-party reviewers into these processes offers scalability, allowing Meta
11 to handle large volumes of flagged content more efficiently. This is particularly important during
12 peak times, such as during major events or crises when the volume of content can surge.

13 191. One of the main challenges with both internal and third-party review is the potential
14 for variations in enforcement. Different reviewers may have different interpretations of policies,
15 leading to inconsistencies. To mitigate this, Meta provides detailed guidelines and training on
16 applying the Community Standards to third-party reviewers.³⁰⁶

17 192. If a human reviewer remains unsure whether a piece of content violates the
18 Community Standards, they can escalate the questionable content to a content policy subject matter
19 expert on the Global Operations team to help the reviewer come to a decision.³⁰⁷ Additionally,
20 where a reviewer is concerned about an imminent credible threat of physical harm, the content is
21 escalated to experts who determine the threat's credibility and whether it should be reported to
22 emergency services.³⁰⁸

23 193. The human review system also provides feedback to the AI system to further

24

³⁰⁶ Meta. (n.d.). *Helping reviewers make the right calls*. Meta Transparency Center. Retrieved
25 June 10, 2025, from <https://transparency.meta.com/enforcement/detecting-violations/making-the-right-calls/>.

26 ³⁰⁷ META3047MDL-001-00000112.

27 ³⁰⁸ META3047MDL-001-00000112.

1 **Highly Confidential (Competitor)**

1 improve it.³⁰⁹ For example, every comment is human reviewed at least twice before it is fed into
2 the AI system to ensure accuracy.³¹⁰ Importantly, when the human review team makes particularly
3 difficult judgment calls, those decisions are fed to the AI system to further mimic them so that
4 more prohibited content can be reviewed proactively.³¹¹

5 **3. Prioritization of Content for Human Review**

6 194. Given the volume of content on Meta's services,³¹² human review on all content
7 reported by users or detected by AI is impracticable. Instead, given the realistic limits of
8 technology and human capacity that Meta faces (as would be the case for any social media service
9 operating at Meta's scale), Meta triages the review of content by humans. The following describes
10 the ranking and review process for organic content (i.e., non-ad content posted by users of
11 Facebook and Instagram).

12 195. Before an item reaches human reviewers, Meta's technology ranks and prioritizes
13 content so that its team of reviewers is presented with and can focus on the most important cases
14 first.³¹³

15 196. Once a piece of content is enqueued for human review, then that piece of content
16 is assigned a predictive integrity value ("pIV"), which is used to determine which piece of content
17
18

19
20 ³⁰⁹ Meta. (2024, Nov. 12). *How Meta prioritizes content for review*. Meta Transparency Center.
<https://transparency.meta.com/policies/improving/prioritizing-content-review/>.

21 ³¹⁰ META3047MDL-003-00036694.

22 ³¹¹ Meta. (2024, Nov. 12). *How enforcement technology works*. Meta Transparency Center.
<https://transparency.meta.com/enforcement/detecting-violations/how-enforcement-technology-works/>.

23
24 ³¹² E.g., Smith, C. (2013, September 18). *Facebook users are uploading 350 million new photos each day*. Business Insider. <https://www.businessinsider.com/facebook-350-million-photos-each-day-2013-9>.

25
26 ³¹³ Meta. (2022, Jan. 19). *How technology helps prioritize review*. Meta Transparency Center.
<https://transparency.meta.com/enforcement/detecting-violations/technology-helps-prioritize-review>.

1 **Highly Confidential (Competitor)**

1 should be reviewed first by the human reviewers.³¹⁴ PIV is calculated based on factors including
2 likelihood of violation, virality, and violation severity, as further explained below.³¹⁵

3 197. At this review stage, I understand a range of signals further contribute to the
4 prediction that a piece of content is policy-violating content for the purpose of prioritization,
5 including Whole Post Integrity Embeddings (“WPIE”). WPIE is a holistic, multimodal, and robust
6 classification layer—in addition to content-specific classifiers described above—that looks at a
7 swath of contextual queues (e.g., image, text, comments) to predict whether a piece of content is
8 policy violating.³¹⁶ WPIE assesses the basic details of the content being detected to determine what
9 kind of violation may be at issue.³¹⁷ As I describe above,³¹⁸ multimodal solutions are quite
10 resource-intensive and require advanced models and large sets of training data to develop and
11 improve upon. At the scale that Meta processes user reports and escalations from first-level
12 classifiers, I find a multimodal solution like that presented by WPIE to be a robust and more-than-
13 reasonable measure toward the protection of users.

14 198. Meta’s review prioritization efforts do not stop there. Meta’s Multi-Armed Bandit
15 (“MAB”) system assesses the content’s severity level and assigns a “predictive severity level”
16 weight to the content based on the policy area it is most likely to potentially violate (i.e., if MAB
17 suspects the content represents a very high severity problem it will get a higher severity weight,

20
21

³¹⁴ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
22 9, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

23 ³¹⁵ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
24 9-10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

25 ³¹⁶ META3047MDL-014-00055649; META3047MDL-014-00073104; *see also*
26 META3047MDL-014-00328102; META3047MDL-072-00001183; META3047MDL-031-
27 00069986.

28 ³¹⁷ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
29 9, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

30 ³¹⁸ See discussion *supra* Part III.B.3.

Highly Confidential (Competitor)

1 while if it is a high or medium severity the assigned weight will be lower).³¹⁹

2 199. Other ranking systems evaluate content virality. The Virality, Estimations,
 3 Prediction, and Applications (“VESPA”) system assesses the piece of content’s virality by
 4 counting each time a user views it for a meaningful amount of time. VESPA assigns a “predicted
 5 virality” value to a content. “Predicted virality value” reflects how widely viewed a particular piece
 6 of content is expected to be.³²⁰ The High Risk Early Review Operations (“HERO”) system is an
 7 end-to-end system used for human review to remove high viral violating content, including content
 8 flagged by VESPA.³²¹

9 200. Based on these evaluations of harm classification, severity, and virality, content is
 10 assigned a pIV to prioritize the piece of content for review.³²² A priority formula then orders
 11 content for review based on pIV.³²³ The highest pIV tasks move to the top of the queue so that the
 12 most viral and potentially harmful content gets reviewed sooner.³²⁴ The queue is designed to be
 13 dynamic, such that it continually brings the highest priority items to the top of the queue for review.
 14 Some content is always prioritized because it is “extremely high severity” content (e.g., content
 15 that may violate Meta’s child safety policies), and even if they have very low virality, the severity
 16 rating will carry it to the top of the review queue.³²⁵

17

18 ³¹⁹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
 19 9, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025). Meta does not assign “weight” to or
 20 otherwise “rank” classifiers. Rather, content is ranked for human review, according to its
 predicted severity. *Id.*

21 ³²⁰ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
 9-10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

22 ³²¹ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
 9-10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

23 ³²² Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

24 ³²³ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

25 ³²⁴ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

26 ³²⁵ Meta Defendants’ Supplemental Responses and Objections to Plaintiffs’ First Interrogatory at
 10, No. 4:22-MD-03047-YGR (N.D. Cal. Feb. 21, 2025).

1 **Highly Confidential (Competitor)**

1 **VI. Effectiveness and Transparency of Meta's Content Moderation Systems**

2 **A. CSERs & the Prevalence Metric**

3 201. Meta has published quarterly Community Standards Enforcement Reports
4 ("CSERs") since 2018.³²⁶ These reports measure prevalence (the estimated percentage of views
5 that were of violating content when compared to all views of content on Facebook or Instagram).³²⁷
6 The most recent CSER data can be found online.³²⁸

7 202. More specifically, prevalence of violating content is estimated using samples of
8 views of content from across Facebook or Instagram.³²⁹ It is calculated as the estimated number
9 of views that showed violating content, divided by the estimated number of total content views on
10 Facebook or Instagram.³³⁰ So, for example, if the prevalence of some category of violating content
11 was 0.05% to 0.06%, that would mean of every 10,000 content views, 5 to 6 on average were of
12 content that violated Meta's standards for adult nudity and sexual activity.³³¹

13 **B. Reasonableness of CSERs for Evaluating the Effectiveness of Content
14 Moderation System**

15 203. Meta's view-based "prevalence" metric is an appropriate measure for evaluating
16 user exposure to harmful content. From a risk-assessment standpoint, what matters is not merely

17
18 ³²⁶ See Rosen, G. (2018, May 15). *Facebook publishes enforcement numbers for the first time*.
19 Meta Newsroom. <https://about.fb.com/news/2018/05/enforcement-numbers/>; (Harwell &
Timberg, 2020).

20 ³²⁷ Meta. (2025, Mar. 6). *Prevalence*. Meta Transparency Center.
21 <https://transparency.meta.com/policies/improving/prevalence-metric/>.

22 ³²⁸ See Meta. (n.d.). *Community standards enforcement report*. Meta Transparency Center.
23 Retrieved June 13, 2025, from <https://transparency.meta.com/reports/community-standards-enforcement/>.

24 ³²⁹ Meta. (2025, Mar. 6). *Prevalence*. Meta Transparency Center.
<https://transparency.meta.com/policies/improving/prevalence-metric/>.

25 ³³⁰ Meta. (2025, Mar. 6). *Prevalence*. Meta Transparency Center.
<https://transparency.meta.com/policies/improving/prevalence-metric/>.

26 ³³¹ Meta. (2025, Mar. 6). *Prevalence*. Meta Transparency Center.
<https://transparency.meta.com/policies/improving/prevalence-metric/>.

Highly Confidential (Competitor)

how much violating content exists on the service, but how much of it is seen. A piece of violating content that is uploaded but quickly removed without exposure poses virtually no harm. Meta's focus on exposure to policy violating content reflects a scientifically grounded approach, consistent with risk-based frameworks used in safety engineering and public health surveillance.³³²

204. The above-described sampling methodology is statistically rigorous and representative, and permits inference about service-wide behavior. This approach reflects current best practices in algorithmic accountability and harm prioritization.

205. Meta also reports several related metrics in its CSERs, including: the number of content items actioned; the “proactive rate” (i.e., the percentage of all content or accounts acted on that Meta found and actioned before users reported them to Meta); the number of appeals; and the number of restored content items. These multiple indicators enable evaluation not only of coverage (recall) but also of precision, fairness, and responsiveness. This multi-metric framework is consistent with the kinds of internal validation and external audit processes used in mature, safety-critical systems.

206. While no measurement system is immune to limitations, Meta's CSER framework is transparent, statistically grounded, and continuously refined. Meta discloses methodology, publishes quarterly data, and permits independent review. In my opinion, this level of transparency and methodological integrity substantially exceeds industry norms.

^{207.} Lastly, it is worth noting that the prevalence metric is well supported in the academic literature.³³³ Leading scholars in the fields of algorithm auditing and service

³³² Sandvig, C., Hamilton, K., Karahalios, K., & Langbort, C. (2014). Auditing algorithms: Research methods for detecting discrimination on internet platforms. *Data and discrimination: converting critical concerns into productive inquiry*, 22(2014), 4349-4357; Binns, R., Van Kleek, M., Veale, M., Lyngs, U., Zhao, J., & Shadbolt, N. (2018, April). 'It's Reducing a Human Being to a Percentage' Perceptions of Justice in Algorithmic Decisions. In *Proceedings of the 2018 CHI conference on human factors in computing systems* (pp. 1-14).

³³³ Yale Law School. (2019, May 23). *Facebook data transparency advisory group releases final report*. <https://law.yale.edu/yls-today/news/facebook-data-transparency-advisory-group-releases-final-report>.

Highly Confidential (Competitor)

1 accountability have advocated for view-based exposure metrics as superior indicators of user harm
2 compared to raw content counts.³³⁴ The CSER framework aligns closely with these
3 recommendations and reflects an empirically responsible posture toward reduction of policy-
4 violating content and transparency:

5 208. First, CSER is based on prevalence, i.e., the estimated proportion of total content
6 views that contained policy-violating material, a measure rooted in audience exposure, not just
7 policy enforcement.

8 209. Second, CSER provides category-specific metrics (e.g., hate speech, graphic
9 violence, adult nudity) and includes confidence intervals, reflecting best practices in measurement
10 transparency and statistical disclosure.

11 210. Third, Meta applies probability-based stratified sampling and then applies manual
12 review procedures, ensuring that prevalence estimates in CSER are both replicable and
13 independently auditable, two key standards recommended by both the Facebook Data
14 Transparency Advisory Group and independent researchers.

15 211. Finally, Meta regularly updates its methodology and discloses changes to its
16 classifiers and enforcement systems over time, which mirrors the kind of iterative evaluation
17 frameworks recommended in transparency and auditing literature.

18 C. Most-Recent CSER Report

19 212. Meta's most recent CSER report was published in February 2025.³³⁵ That report
20 included measurements on prevalence, content actioned, proactive rate, appealed content, and
21 restored content, spanning the policy areas of violence and criminal behavior, safety, objectionable

22 _____
23 ³³⁴ Gorwa, R., Binns, R., & Katzenbach, C. (2020). Algorithmic content moderation: Technical
24 and political challenges in the automation of platform governance. *Big Data & Society*, 7(1),
2053951719897945; Ribeiro, M. H., Ottoni, R., West, R., Almeida, V. A., & Meira Jr, W. (2020,
25 January). Auditing radicalization pathways on YouTube. In Proceedings of the 2020 conference
on fairness, accountability, and transparency (pp. 131-141).

26 ³³⁵ See Meta. (n.d.). *Community standards enforcement report*. Meta Transparency Center.
27 Retrieved June 13, 2025, from <https://transparency.meta.com/reports/community-standards-enforcement/>.

Highly Confidential (Competitor)

1 content, and integrity and authenticity. Meta presently defines these metrics as follows:

2 213. Prevalence: estimated percentage of views that were of violating content (defined
3 in more detail above).³³⁶

4 214. Content actioned: the number of pieces of content (such as posts, photos, videos or
5 comments) or accounts taken action on for going against standards; shows the scale of enforcement
6 activity.³³⁷

7 215. Proactive rate: the percentage of all content or accounts acted on that Meta found
8 and actioned *before* users reported them to Meta; used as an indicator of how effectively Meta
9 detects violations.³³⁸

10 216.Appealed content: the number of pieces of content (such as posts, photos, videos or
11 comments) that people appeal after Meta takes action on for going against policies.³³⁹

12 217. Restored content: the number of pieces of content (such as posts, photos, videos or
13 comments) restored after Meta originally took action on them.³⁴⁰

14 218. As of the Q1 2025 report, CSER published prevalence metrics from Instagram and
15 Facebook for the following relevant categories of content:

16 *Table 6. Adult Nudity & Sexual Content CSER*

Metric	Instagram	Facebook
Prevalence	Lower bound and upper bound: 0.04%	Lower bound: 0.05%; Upper bound: 0.06%
Content Actioned	9.8M	64.7M
ContentAppealed	1.6M	2.5M

20 ³³⁶ Meta. (2025, Mar. 6). *Prevalence*. Meta Transparency Center.

21 <https://transparency.meta.com/policies/improving/prevalence-metric/>.

22 ³³⁷ Meta. (2023, Nov. 7). *Content actioned*. Meta Transparency Center.

23 <https://transparency.meta.com/policies/improving/content-actioned-metric/>.

24 ³³⁸ Meta. (2023, Feb. 22). *Proactive rate*. Meta Transparency Center.

25 <https://transparency.meta.com/policies/improving/proactive-rate-metric/>.

26 ³³⁹ Meta. (2022, Nov. 18). *Appealed content*. Meta Transparency Center.

27 <https://transparency.meta.com/policies/improving/appealed-content-metric/>.

28 ³⁴⁰ Meta (2022, Oct. 4). *Restored content*. Meta Transparency Center.

<https://transparency.meta.com/policies/improving/restored-content-metric/>.

Highly Confidential (Competitor)

Metric	Instagram	Facebook
Proactive Rate	97.7%	94.9%
Restored Content	825.1K with appeal; 97.8K without appeal	1.1M with appeal; 191.2K without appeal

Table 7. Bullying & Harassment CSER

Metric	Instagram	Facebook
Prevalence	Lower bound: 0.05%; Upper bound: 0.06%	Lower bound: 0.07%; Upper bound: 0.08%
Content Actioned	5.2M	5.1M
ContentAppealed	778.3K	904K
Proactive Rate	88%	73.3%
Restored Content	196.9K with appeal; 79.4K without appeal	143.6K with appeal; 23.4K without appeal

Table 8. Child Endangerment (nudity and physical abuse) CSER

Metric	Instagram	Facebook
Prevalence	N/A	N/A
Content Actioned	616K	1.5M
ContentAppealed	58.8K	163K
Proactive Rate	98%	97.2%
Restored Content	23.5K with appeal; 7.2K without appeal	32KK with appeal; 42.7K without appeal

Table 9. Child Endangerment (sexual exploitation) CSER

Metric	Instagram	Facebook
Prevalence	N/A	N/A
Content Actioned	1.5M	4.6M
ContentAppealed	93.2K	398.1K
Proactive Rate	94.5%	95.8%
Restored Content	28.9K with appeal; 32.5K without appeal	127K with appeal; 287.4K without appeal

Table 10. Suicide & Self-Injury CSER

Metric	Instagram	Facebook
Prevalence	Upper bound: 0.05%	Upper bound: 0.05%
Content Actioned	9.9M	6.8M
ContentAppealed	118.3K	231.6K
Proactive Rate	99.10%	98.9%
Restored Content	24.7K with appeal; 1.2K without appeal	41.3K with appeal; 1.5K without appeal

Highly Confidential (Competitor)1 *Table 11. Violent & Graphic Content CSER*
2

Metric	Instagram	Facebook
Prevalence	Lower bound and Upper bound: 0.06%	Lower bound: 0.09%; Upper bound: 0.09%
Content Actioned	4.9M	10.7M
ContentAppealed	26K	64.2K
Proactive Rate	97.5%	98.0%
Restored Content	10.3K with appeal; 31.4K without appeal	12.3K with appeal; 12.5K without appeal

7 219. Meta's most recent Community Standards Enforcement Report discloses several
8 quantitative metrics that collectively provide a robust and multi-dimensional picture of the
9 service's content moderation efficacy.

10 220. Proactive Rate indicates the percentage of violating content detected and actioned
11 before users report it. High proactive rates, such as the 98.9% proactive detection of suicide and
12 self-injury content on both Facebook and Instagram, demonstrate the success of Meta's AI-driven
13 classifier systems in identifying policy-violating content without relying on user intervention.

14 221. The fact that only a relatively small portion of actioned content is appealed, and
15 that an even smaller subset of those appeals result in restoration after human review, suggests that
16 the enforcement system is both accurate and subject to meaningful user recourse.
17 In conclusion, while no moderation system is perfect, Meta's publication of granular, view-based
18 transparency and enforcement norms across the social media services industry. These features lend
19 empirical and comparative support to the reasonableness and good-faith design of its moderation
20 framework.

21 **D. Expert Assessment of the Algorithmic Stress Test of Instagram's Reels Surface**

22 222. In this section, I evaluate the purported testing protocol of Instagram accounts that
23 Mr. Arturo Bejar, a former employee and contingent worker at Meta who worked on user safety
24 and security, claims he conducted from late 2023 to early 2024 "to see what the experience was of
25

Highly Confidential (Competitor)

1 a teenager on Instagram.”³⁴¹ To evaluate Mr. Bejar’s testing, I have reviewed relevant portions of
 2 his deposition transcripts and the accompanying exhibits, with particular attention to Exhibit 12.³⁴²
 3 My opinion of Mr. Bejar’s testing protocol can be summarized as follows:

4 223. *First*, although Mr. Bejar does not characterize his testing protocol as such,
 5 functionally, Mr. Bejar appears to have conducted a truncated form of algorithmic stress test on
 6 Instagram’s Reels surface.

7 224. *Second*, as an algorithmic stress test, Mr. Bejar’s protocol and methodology suffer
 8 from fatal deficiencies that render the results scientifically unreliable.

9 225. I reach these conclusions based on the following evaluation.

10 **1. Mr. Bejar Conducted a Truncated Algorithmic Stress Test**

11 226. Mr. Bejar claims he evaluated “a few test accounts” “to see what the experience
 12 was of a teenager on Instagram.”³⁴³ Though Mr. Bejar does not identify his evaluation as such,
 13 Mr. Bejar’s deposition testimony and accompanying exhibits indicate that he attempted an
 14 algorithmic stress test of Instagram’s Reels content recommendation algorithm.

15 227. Algorithmic stress tests are forms of audits that are performed on algorithmic
 16 systems to test and replicate the behavior of the system under certain controlled conditions. When
 17 properly conducted, stress tests are part of an established toolkit of algorithmic auditing techniques
 18 that are acknowledged as valid inquiry tools by the scientific community.³⁴⁴

19
 20 ³⁴¹ Arturo Bejar Dep. Vol. 1 at 197:1-4, Apr. 7, 2025.

21 ³⁴² See Arturo Bejar Dep. Vol. 1, Apr. 7, 2025; Arturo Bejar Dep. Vol. 2, Apr. 8, 2025; Arturo
 22 Bejar Dep. Vol. 3, Apr. 9, 2025.

23 ³⁴³ Arturo Bejar Dep. Vol. 1 at 197:1-4, Apr. 7, 2025.

24 ³⁴⁴ See the vast body of literature that rigorously outlines protocols for algorithmic audits,
 25 including:

26 (1) Sandvig, C., Hamilton, K., Karahalios, K., & Langbort, C. (2014). *Auditing algorithms:*
 27 *Research methods for detecting discrimination on internet platforms*. Data and
 28 discrimination: converting critical concerns into productive inquiry, 22(2014), 4349-
 4357.

Highly Confidential (Competitor)

228. Though his methods were scientifically flawed, Mr. Bejar appears to have attempted an assessment of the behavior of Instagram's Reels algorithm under particular conditions, which suggests he attempted to conduct an algorithmic stress test. Using an Apple iPhone 11 and an iPad at indeterminate times, Mr. Bejar created a small yet indeterminate number of synthetic Instagram accounts through unique iCloud email addresses. Four accounts purported to represent minor female users, an indeterminate number of accounts purported to represent minor male users, and one account purported to represent a 23- or 25-year-old of indeterminate sex and gender, as a purported control.³⁴⁵ Mr. Bejar manually steered some of these accounts toward content categories on Reels he deemed harmful or inappropriate (e.g., sexually suggestive material, self-harm-related content), or media involving body image issues or minors in potentially suggestive contexts. These interactions were recorded via screenshots and screen-capture videos, which were then presented as purported evidence that Instagram's recommender system promotes or amplifies harmful content to underage users.

229. While the concept of a “stress test” has a recognized place in both academic and engineering practice, Mr. Bejar’s implementation of this method fails to conform to any scientifically accepted protocol for conducting such audits. In the context of algorithmic systems, a legitimate stress test is a *controlled, repeatable, and statistically grounded procedure* that deliberately subjects a system to adverse or edge-case inputs to evaluate its robustness, performance limits, or failure modes. These tests must be carefully designed to isolate causal relationships between input behaviors and system responses and must be accompanied by appropriate control conditions, transparency of design, and rigorous statistical evaluation. Mr.

(2) Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T., Hutchinson, B., ... & Barnes, P. (2020, January). *Closing the AI accountability gap: Defining an end-to-end framework for internal algorithmic auditing*. In Proceedings of the 2020 conference on fairness, accountability, and transparency (pp. 33-44).

³⁴⁵ BEJAR0002658 (indicating the account purported to represent a 23-year-old); Arturo Bejar Dep. Vol. 2 at 424:16-19, Apr. 8, 2025 (indicating the account purported to represent a 25-year-old).

1 **Highly Confidential (Competitor)**

1 Bejar's so-called test contains none of these critical elements.

2 230. If the goal is to evaluate how recommendation systems operate under stress or
3 boundary conditions, particularly in safeguarding minors, a scientifically valid audit would
4 ordinarily follow a transparent, causally informative experimental design. This entails several
5 methodological components:

6 231. ***Clearly Defined Harm Categories:*** First, researchers must pre-register³⁴⁶
7 operational definitions for the types of harmful content being audited (e.g., pro-eating disorder
8 imagery, sexually suggestive material involving minors, graphic self-harm). These definitions
9 should be aligned with platform policies and third-party standards (e.g., WHO, NSPCC) to ensure
10 validity.

11 232. ***Standardized Behavioral Profiles:*** Test accounts must simulate real user behaviors
12 (e.g., pausing on sensitive content, following triggering accounts, liking specific hashtags) using
13 pre-scripted behavioral templates. These should be uniformly applied across conditions to reduce
14 confounds.

15 233. ***Randomized Assignment and Control Groups:*** A core principle of causal inference
16 is randomization. Accounts should be randomly assigned to different behavioral archetypes (e.g.,
17 neutral, vulnerable, exploratory), with control accounts that exhibit neutral or passive browsing.
18 This creates counterfactual conditions that allow for attributing observed effects to the manipulated
19 behaviors.

20 234. ***Temporal Controls:*** Since recommendation algorithms evolve over time,

21 ³⁴⁶ In simplistic terms, *pre-registration* means that before starting a study or audit, the
22 researcher(s) would state and pre-register (write down in advance) what phenomenon or effect
23 they're going to look for, how they'll look for it, and how they'll measure it. In the specific case,
24 this would include clearly defining what counts as harmful content (like what exactly qualifies as
25 pro-eating disorder imagery or sexually suggestive material involving minors). Pre-registration
26 of hypotheses and methodologies are scientific best practices ensuring that (i) researchers don't
27 change the rules halfway through to get the results they want; (ii) The employed process is
transparent and replicable, so others can check the work; (iii) This reduces bias by committing to
a plan before seeing any results. As a common practice in high-quality scientific studies, pre-
registration promotes honesty, objectivity, and accountability, especially when the findings may
be used to inform policy or public decisions.

19 **Highly Confidential (Competitor)**

1 experiments must be run in parallel and repeated over time to account for drift, retraining cycles,
2 or time-of-day effects.

3 235. ***Comprehensive Logging:*** Every piece of recommended content should be captured
4 with timestamped metadata including its position in the feed, the algorithmic pathway (e.g., For
5 You vs. Search), source account characteristics, and user engagement logs. This dataset should
6 then be made available for auditing and replication purposes.

7 236. ***Blinded Annotation:*** To evaluate exposure, all logged content should be labeled by
8 multiple blinded raters using a predefined codebook. Reliability statistics (e.g., Cohen's κ or
9 Krippendorff's α) should then be reported to ensure consistency.

10 237. ***Statistical Analysis:*** Analysis must use appropriate hypothesis testing and effect
11 size estimation. For binary outcomes (e.g., exposure to harmful content), chi-squared tests or
12 logistic regression may be appropriate. For prevalence comparisons, z-tests for proportions or
13 bootstrapped confidence intervals can quantify uncertainty. Causal inference techniques such as
14 difference-in-differences (DiD) or matched group comparisons may also be employed when
15 longitudinal or quasi-experimental designs are available.

16 238. By adhering to these principles, researchers can credibly isolate causal effects,
17 quantify the reliability and scale of potential exposure harms, and ensure that their methodology
18 is both scientifically transparent and statistically valid.

19 **2. Mr. Bejar's Algorithmic Stress Test Suffers from Fatal
20 Methodological Flaws Rendering the Results Scientifically Unreliable**

21 239. Mr. Bejar's methodology suffers from the following core deficiencies:

22 240. ***Lack of Experimental Design:*** There is no evidence that Mr. Bejar employed an
23 experimental design worthy of the name. A valid experimental design necessarily includes a
24 predefined hypothesis, experimental or control groups, randomized treatments, and standardized
25 account initialization protocols. Mr. Bejar's test lacks evidence of each of these features. His
26 experimental approach appears arbitrary in critical respects. He conceived of the idea
27 independently after reading an article in the Wall Street Journal about the prevalence of certain

Highly Confidential (Competitor)

1 content on Instagram,³⁴⁷ but he does not appear to have consulted any sources or collaborated with
2 any qualified experts to aid in the design or direction of the experiment.³⁴⁸ Further, there is no
3 indication Mr. Bejar developed a set of specific criteria for the type of content recommendations
4 he sought to assess, identified the relevant user populations, or otherwise set any particularized,
5 time-limited goals for the test. On the contrary, Mr. Bejar took a casual approach guided by
6 convenience and his personal whims. Mr. Bejar created “approximately seven” synthetic accounts
7 as minor males and females³⁴⁹ and conducted “initial testing” (which lasted for several minutes at
8 most for each account) to see whether Reels would recommend “racy” or “violent” content.³⁵⁰ But
9 at some point he stopped checking most of the accounts after initial testing because he found the
10 search results emotionally distressing.³⁵¹ During this pause, without methodological explanation,
11 Mr. Bejar occasionally searched the phrase “I want to hurt myself” on one or more of the
12 accounts,³⁵² and may have resumed other searches, but would stop soon again because he found
13 the content disturbing.³⁵³ When Mr. Bejar returned more fully to his project approximately a year
14 after he started, he no longer checked the minor male accounts and instead focused on the minor
15 female accounts because he “had all of them set up in a single phone and that made it easier.”³⁵⁴
16 In short, Mr. Bejar allowed his personal sensitivities and whims delineate his execution of the test.
17 In academic literature, robust algorithmic audits require pre-registration or structured protocols
18 that specify how synthetic accounts are configured, how engagement behaviors are scripted, and
19 what constitutes a treatment versus a control condition. Mr. Bejar provides no such evidence or

20 ³⁴⁷ Arturo Bejar Dep. Vol. 3 at 995:13-25, Apr. 9, 2025.

21 ³⁴⁸ Arturo Bejar Dep. Vol. 3 at 1008-09, Apr. 9, 2025.

22 ³⁴⁹ Arturo Bejar Dep. Vol. 3 at 994:8-9, Apr. 9, 2025.

23 ³⁵⁰ Arturo Bejar Dep. Vol. 1 at 199:8-15, Apr. 7, 2025 (explaining that initial testing lasted
anywhere from “8 to 12 minutes, sometimes 18”).

24 ³⁵¹ Arturo Bejar Dep. Vol. 3 at 998:8-11, Apr. 9, 2025.

25 ³⁵² Arturo Bejar Dep. Vol. 3 at 998:18-21, Apr. 9, 2025

26 ³⁵³ Arturo Bejar Dep. Vol. 3 at 999:5-8, Apr. 9, 2025

27 ³⁵⁴ Arturo Bejar Dep. Vol. 3 at 999:1-8, 997:12-17, Apr. 9, 2025.

Highly Confidential (Competitor)

documentation. The absence of these elements renders the protocol scientifically uninterpretable and prevents any possibility of causal inference.

241. ***Confirmation and Interaction Bias:*** The test accounts were clearly guided through intentional behaviors designed to surface borderline or policy-violating content. This includes searching for sensitive hashtags and engaging with content flagged as suggestive. For example, on one test account, Mr. Behar expressly “searched for some gymnastics content” which eventually led him to “videos of very young girls in different kinds of outfits.”³⁵⁵ On other accounts, Mr. Bejar searched the phrase “I want to hurt myself” an indeterminate number of times over a yearlong period.³⁵⁶ On other test accounts, Mr. Bejar claims he did not conduct any searches, but quickly swiped through Reels videos without finishing until he came across content he deemed harmful (e.g., “a little violent” or “a little racy”).³⁵⁷ If the user posting the video appeared to Mr. Bejar to be under the age of thirteen, Mr. Bejar would watch the video to completion.³⁵⁸ These forms of active user engagement are known to influence recommender systems,³⁵⁹ which are designed to tailor content based on user behavior. This creates a feedback loop,³⁶⁰ known in the literature as *interaction bias*, which reinforces the types of content a user engages with. Without neutral baselines or randomized interaction behaviors, the experiment cannot distinguish between service-driven recommendations and user-driven discovery.³⁶¹

242. *Opacity and Poor Documentation*: The documentation accompanying Mr. Bejar's

³⁵⁵ Arturo Bejar Dep. Vol. 1 at 198:20-25, Apr. 7, 2025.

³⁵⁶ Arturo Bejar Dep. Vol. 2 at 402:21-24, Apr. 8, 2025.

³⁵⁷ Arturo Bejar Dep. Vol. 1 at 199:1-15, Apr. 7, 2025.

³⁵⁸ Arturo Bejar Dep. Vol. 1 at 199:10-12, Apr. 7, 2025

³⁵⁹ Baeza-Yates, R. (2020, September). Bias in search and recommender systems. In Proceedings of the 14th ACM conference on recommender systems (pp. 2-2).

³⁶⁰ Mansouri, M., Abdollahpouri, H., Pechenizkiy, M., Mobasher, B., & Burke, R. (2020, October). *Feedback loop and bias amplification in recommender systems*. In Proceedings of the 29th ACM international conference on information & knowledge management (pp. 2145-2148).

³⁶¹ Khenissi, S., & Nasraoui, O. (2020). Modeling and counteracting exposure bias in recommender systems. arXiv preprint arXiv:2001.04832.

Highly Confidential (Competitor)

1 materials consists primarily of editorialized file names, screen recordings, and fragmented screen
2 captures. There is no formalized protocol outlining how many accounts would be used, over what
3 time period the testing would be conducted, whether any efforts would be made to replicate
4 findings, or how the outputs would be categorized and labeled. The only written record of the test
5 consists of the so-called “test protocol” that Mr. Bejar drafted in 2025, long after the start of Mr.
6 Bejar’s inquiry in 2023.³⁶² The fact that the protocol was drafted long after the start of testing
7 provides little confidence that the results of Mr. Bejar’s efforts followed from the guidance the
8 protocol sets forth. Moreover, Mr. Bejar’s protocol is high-level, conclusory, and lacks all of the
9 details outlined above that would even permit drawing reliable connections between the protocol
10 and the results. Without purporting to identify all examples, the following are several of the critical
11 gaps between Mr. Bejar’s claims and his drafted protocol. Most critically, the protocol offers no
12 transparency around whether content was organically recommended by the service or was the
13 result of user searches or manual navigation. Further, the protocol offers no support for Mr. Bejar’s
14 assertion that he recorded video from each testing session,³⁶³ never mind his admission that in
15 some cases he only recorded part of the testing session.³⁶⁴ Further, the protocol makes no mention
16 that Mr. Bejar paused testing nor that he abandoned his consideration of the minor male accounts
17 upon resuming his experiment.³⁶⁵ Nor does the protocol support Mr. Bejar’s claim that he planned
18 testing sessions in advance, nor that Mr. Bejar would erase the iPhone 11 and redownload
19 Instagram between testing sessions for each synthetic account.³⁶⁶ As another example, Mr. Bejar
20 claims he used, at various times, an iPhone 11, and at other times an iPad, but the protocol mentions
21

22³⁶² Arturo Bejar Dep. Vol. 3 at 1001:12-16, Apr. 9, 2025; *see BEJAR0002658.*

23³⁶³ Arturo Bejar Dep. Vol. 3 at 1020:9-22, Apr. 9, 2025.

24³⁶⁴ Arturo Bejar Dep. Vol. 3 at 1021:1-5, Apr. 9, 2025. For example, Mr. Bejar admitted that he
25 did not consistently capture some of his testing sessions wherein he searched the phrase “I want
to hurt myself” and similar phrases. *Id.* at 1027-28.

26³⁶⁵ Arturo Bejar Dep. Vol. 3 at 1005:1-18, Apr. 9, 2025.

27³⁶⁶ Arturo Bejar Dep. Vol. 3 at 1019, Apr. 9, 2025.

1 **Highly Confidential (Competitor)**

2 only the former.³⁶⁷ This lack of documentation makes the protocol entirely non-replicable and non-auditable.

3 243. **No Statistical Rigor or Validity:** Mr. Bejar provides no statistical measurement of
4 exposure rates, prevalence of harmful content, latency of content surfacing, or control comparisons
5 across different account behaviors. He does not define any metrics for measuring algorithmic
6 performance (e.g., precision, recall, rate of harmful content exposure), nor does he employ any
7 tools for validating whether the content exposure observed was anomalous or consistent with
8 typical service behavior. The sample size is undefined or too small to yield a statistically
9 appropriately powered study; no effort is made to perform statistical analysis across multiple trials.
10 The results that are shown are the byproduct of “statistical cherry-picking”,³⁶⁸ i.e., the malpractice
11 of selecting evidence that corroborates a position, while discarding evidence of the contrary: Mr
12 Bejar excluded the results from a subset of synthetic accounts he created, approximately half of
13 the total, without any explanation or justification. These deficiencies disqualify the methodology
14 from being considered scientifically valid.

15 244. **Failure to Disentangle Algorithmic Output from User Behavior:** At its core, Mr.
16 Bejar’s approach confounds the effects of service recommendation mechanisms with user steering
17 behavior.³⁶⁹ Many of the outputs catalogued appear to follow naturally from the actions taken by
18 the user (e.g., lingering on, liking, or replaying suggestive videos while ignoring any other type of
19 content). Without controlling for these behaviors, or at least holding them constant across
20 experimental conditions, it is impossible to determine whether the recommender system is
21 surfacing harmful content organically or simply responding to input behavior, as designed. No

22
23 ³⁶⁷ Arturo Bejar Dep. Vol. 3 at 1011-13, Apr. 9, 2025.

24 2368 *Cherry picking*. (n.d.). Wikipedia. Retrieved June 13, 2025, from
25 https://en.wikipedia.org/wiki/Cherry_picking#:~:text=Cherry%20picking%2C%20suppressing%20evidence%2C%20or,that%20may%20contradict%20that%20position.

26 269 Chaney, A. J., Stewart, B. M., & Engelhardt, B. E. (2018, September). *How algorithmic
27 confounding in recommendation systems increases homogeneity and decreases utility*. In
Proceedings of the 12th ACM conference on recommender systems (pp. 224-232).

1 **Highly Confidential (Competitor)**

1 counterfactuals are presented to test alternate explanations.

2 245. Accordingly, Mr. Bejar's findings cannot be considered reliable or generalizable
3 assessments of how Instagram's recommendation systems function. His results are better
4 characterized as anecdotal observations lacking the scientific safeguards required to support
5 meaningful conclusions. Simply put, what Mr. Bejar conducted is not a "stress test" in any
6 scientifically accepted sense; it is a series of anecdotal observations devoid of methodological
7 grounding.

8 246. Mr. Bejar's findings do not provide meaningful evidence of how Instagram's
9 recommender systems function either under normal circumstances or under simulated stress
10 conditions. They do not demonstrate systemic failure in content moderation, algorithmic bias, or
11 harmful amplification. They provide no insight into algorithmic thresholds, classifier behavior, or
12 enforcement mechanisms. In short, Mr. Bejar's materials are not only methodologically unsound,
13 but they are also fundamentally unreliable and should carry no evidentiary weight in any serious
14 inquiry into service integrity.

15 **VII. Meta's Content Moderation Policies & Enforcement Align with Industry Best
16 Practices**

17 247. Meta's content moderation policies, the technology and review processes that
18 enforce them, as well as the public reporting of those efforts align with best practices in the
19 technology industry.

20 248. The Digital Trust & Safety Partnership ("DTSP") has developed a widely accepted
21 framework that articulates industry best practices for managing content- and conduct-related risks
22 across five key commitments: (i) risk evaluation in product development, (ii) product governance,
23 (iii) enforcement operations, (iv) iterative improvement, and (v) transparency.³⁷⁰ In my expert
24 opinion, Meta's content moderation infrastructure and governance systems not only align with the

25
26 ³⁷⁰ Digital Trust & Safety Partnership. (n.d.). *Trust & Safety Best Practices Framework*.
27 Retrieved June 11, 2025, from https://dtspartnership.org/wp-content/uploads/2021/04/DTSP_Best_Practices.pdf.

Highly Confidential (Competitor)

1 DTSP's framework but in many areas serve as a leading implementation of these best practices.

2 **A. Product Development: Proactive Risk Identification and Mitigation**

3 249. Meta integrates Trust & Safety considerations from the earliest stages of product
4 development. As detailed in this report, Meta convenes biweekly Product Policy Forums that
5 include safety and cybersecurity policy teams, product managers, and legal and policy leads to
6 assess risk before feature deployment. The company's content policy team, informed by internal
7 analyses and external consultations (e.g., with mental health and child safety organizations),
8 participates directly in product iteration. This mirrors DTSP's recommendation that companies
9 embed risk assessment, stakeholder input, and feedback loops into the development lifecycle.

10 **B. Product Governance: Transparent and Evolving Policy Structures**

11 250. Meta maintains a layered governance structure consisting of its public Community
12 Standards, internal Policy Labs, and Borderline Content Policies. These are developed with input
13 from academics, civil society groups, and domain experts. Changes to policy are logged and made
14 publicly visible. These practices reflect DTSP's emphasis on explainable rulemaking, user-facing
15 clarity, community input, and documented internal interpretative standards.

16 **C. Enforcement Operations: Scalable and Responsive Moderation Infrastructure**

17 251. Meta operationalizes product governance through a hybrid enforcement system
18 involving (i) scalable AI classifiers and proactive detection, (ii) global human review teams
19 including specialized "ring-fenced" experts, (iii) tiered prioritization mechanisms such as
20 Predictive Integrity Value (pIV) and the HERO system, and (iv) robust user reporting and appeals
21 systems. Meta also provides employee wellness support and workload safeguards for reviewers,
22 directly aligning with DTSP guidance on workforce resilience, reporting tools, and scaled
23 enforcement workflows.

24 **D. Iterative Improvement: Data-Driven Learning and Policy Adjustment**

25 252. Meta continuously updates its classifiers, policy enforcement thresholds, and risk
26 mitigation strategies based on feedback from audits, real-world incidents, and statistical
27 performance reviews. The use of systems like Label Accuracy Management (LAM), the

Highly Confidential (Competitor)

1 enforcement metrics reported in CSERs (e.g., proactive detection rates), and ongoing adjustments
2 based on expert consultation exemplify DTSP's call for structured learning loops and explainable
3 frameworks for policy evolution.

4 **E. Transparency and Accountability: Industry-Leading Disclosures**

5 253. Meta publishes detailed Community Standards Enforcement Reports (CSERs)
6 quarterly, providing metrics on prevalence, proactive rate, appeals, and restoration. These reports
7 include confidence intervals, describe classifier methodology, and enable public scrutiny. Meta
8 also engages with researchers through dataset access initiatives, sponsors benchmarking
9 challenges (e.g., Deepfake Detection and Hateful Memes), and incorporates in-product
10 enforcement signals such as warning screens and user notices. These practices are consistent with
11 DTSP's call for public-facing policy disclosures, support for academic collaboration, and periodic
12 reporting to stakeholders.

13 254. In sum, Meta not only fulfills the commitments described in the DTSP Best
14 Practices Framework but does so at a scale and level of transparency that sets a benchmark for the
15 industry. Based on my review of Meta's content moderation infrastructure, I conclude that Meta
16 is a paradigmatic example of a DTSP "Practicing Company."

17 **VIII. Conclusions**

18 255. Based on my review of the evidence, internal documentation, publicly disclosed
19 data, and relevant technical and scientific literature, I offer the following expert findings regarding
20 Meta's content moderation systems:

21 256. Meta has implemented one of the most technically advanced and operationally
22 mature content moderation infrastructures in the industry. Its multi-pronged approach includes AI-
23 based classifiers and other automated detection systems, human review pipelines, escalation
24 procedures, and reactive user-report channels, all of which are integrated into a coordinated
25 enforcement framework.

26 257. Meta's Community Standards Enforcement Reports ("CSERs") provide a
27 scientifically sound and statistically robust mechanism for evaluating the effectiveness of content

Highly Confidential (Competitor)

1 moderation. The “prevalence” metric, in particular, offers a meaningful proxy for potential user
2 harm by measuring actual exposure to violating content.

3 258. Meta’s moderation systems prioritize high-severity harms, such as suicide, self-
4 injury, and child exploitation, and demonstrate exceptionally high proactive detection rates in these
5 categories, often exceeding 98%. These rates reflect the maturity and effectiveness of Meta’s AI
6 classifiers, as well as the prioritization mechanisms embedded in its enforcement architecture.

7 259. The company has demonstrated a transparent and iterative approach to measuring,
8 reporting, and refining its moderation systems. The inclusion of appeal and restoration metrics,
9 alongside published methodologies, underscores a commitment to accountability and continuous
10 improvement.

11 260. Meta’s systems are not static but are designed to evolve in response to emerging
12 harms, adversarial tactics, cultural considerations, and regulatory constraints. Classifier thresholds,
13 enforcement priorities, and borderline content protocols are continuously updated based on new
14 evidence and service dynamics.

15 261. Allegations that Meta’s recommendation systems promote harmful content,
16 including those based on anecdotal stress tests such as the one conducted by Mr. Bejar, are not
17 grounded in scientifically valid methodologies. Such audits fail to meet basic requirements for
18 control, replication, statistical rigor, or causal inference, and should not be afforded evidentiary
19 weight in serious evaluations of service integrity.

20 262. Taken together, the evidence supports the conclusion that Meta’s content
21 moderation systems are robust, reasonably designed, and continuously improving in ways that are
22 consistent with both academic best practices and industry standards.

23
24 
25

Emilio Ferrara, Ph.D.

26 June 13, 2025
27
28

*Highly Confidential (Competitor)*1 Appendix A: Table of Harm-Specific Classifiers
2

Plain-language classifier name	Alpha-numeric code	Launch date of current iteration	Threshold value for autodeletion	Threshold value for soft actions
Suicide & Self-Injury				
SSI Classifier	F590828509	08/12/2024	.94	<p>Facebook: various enforcements >P50</p> <ul style="list-style-type: none"> • EU Demotion on Newsfeed <ul style="list-style-type: none"> ◦ Suicide: 0.8 ◦ Self-injury: 0.8 • EU Deboost on Newsfeed <ul style="list-style-type: none"> ◦ Suicide: 0.3 ◦ Self-injury: 0.4 • Non-EU Demotion on Newsfeed <ul style="list-style-type: none"> ◦ Suicide: 0.5 ◦ Self-injury: 0.5 • Non-EU Deboost on Newsfeed <ul style="list-style-type: none"> ◦ Suicide: 0.3 ◦ Self-injury: 0.4 • Age-based Demotion on Newsfeed <ul style="list-style-type: none"> ◦ Suicide: 0.3 ◦ Self-injury: 0.3 • Age-based Deboost on Newsfeed <ul style="list-style-type: none"> ◦ Suicide: 0.05 ◦ Self-injury: 0.05 • Age gating on Newsfeed <ul style="list-style-type: none"> ◦ Suicide: 0.1 ◦ Self-injury: 0.1 • Non-rec filtering <ul style="list-style-type: none"> ◦ Suicide: 0.105 ◦ Self-injury: 0.25 • Filtering demotion on Stories <ul style="list-style-type: none"> ◦ Self-injury: 0.25 • Age gating on Stories <ul style="list-style-type: none"> ◦ Self-injury: 0.25 • Age gating on Stories (reduce more) <ul style="list-style-type: none"> ◦ Self-injury: 0.1

Highly Confidential (Competitor)

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28

Instagram: various enforcements > P40

- Age gating
 - Suicide: 0.89
 - Self-injury: 0.71
- Age gating High COE
 - Suicide: 0.6
 - Self-injury: 0.24
- Non-rec filtering
 - Suicide: 0.6
 - Self-injury: 0.24
- Hashtag filtering
 - Suicide: 0.80
 - Self-injury: 0.80
- Media SERP
 - Suicide: 0.85
 - Self-injury: 0.85
- Age-filtering
 - Suicide: 0.6
 - Self-injury: 0.24
- Demotion Feed teen non-EU
 - Suicide: 0.7
 - Self-injury: 0.7
- Demotion Feed teen EU
 - Suicide: 0.85
 - Self-injury: 0.85
- Demotion Stories teen
 - Suicide: 0.7Self-injury: 0.7

SSI Borderline Classifier f551409441 04/15/2024 N/A Facebook: various enforcements above >P10

- Age gating on Newsfeed: 0.1
- P-non MSI: 0.3
- Non-rec filtering: 0.3
- Age gating on Stories: 0.75
- Search demotion: 0.5

Instagram: various enforcements > P40

- Age gating: 0.91
- Age gating High CoE: 0.54
- Non-rec filtering: 0.34
- Media SERP: 0.85

Highly Confidential (Competitor)

1				<ul style="list-style-type: none"> • Hashtag filtering: 0.80 • Age filtering: 0.25
Eating Disorder				
2	ED Classifier	f638461582	09/01/2024	N/A
3	4	5	6	Facebook
7	8	9	10	<ul style="list-style-type: none"> • EU Demotion on Newsfeed: 0.8 • EU Deboost on Newsfeed: 0.15 • Non-EU Demotion on Newsfeed: 0.5 • Non-EU Deboost on Newsfeed: 0.15 • Age based Demotion on Newsfeed: 0.3 • Age based Deboost on Newsfeed: 0.05 • Age gating on Newsfeed: 0.1 • Non-rec filtering: 0.061
11	12	13	14	Instagram: various enforcements >P40
15	16	17	18	<ul style="list-style-type: none"> • Non-rec filtering: 0.81 • Media SERP: 0.85 • Hashtag filtering: 0.80 • Demotion Feed teen non-EU: 0.7 • Demotion Feed teen EU: 0.85 • Demotion Stories teen: 0.7
19	ED Borderline Classifier	f551542423	04/15/2024	N/A
20	21	22	23	Facebook
24	25	26	27	<ul style="list-style-type: none"> • Age gating on Newsfeed: 0.1 • P-non MSI: 0.3 • Non-rec filtering: 0.33 • Search demotion: 0.5
28	Bullying & Harassment			

Highly Confidential (Competitor)

1	FB Comment	8165225660248547	10/20/2024	Bullying and Harassment = 0.95 Borderline Hostile Speech = n/a	Bullying and Harassment: <ul style="list-style-type: none">Filtering on Comments: 0.432Age filtering comments: 0.275 Borderline Hostile Speech: <ul style="list-style-type: none">Filtering on Comments: 0.6Age filtering on Comments: 0.2
7	FB Post	8655505921181700	10/23/2024	Bullying and Harassment = 0.95 Borderline Hostile Speech = n/a	Bullying and Harassment: <ul style="list-style-type: none">Age-based Non-rec Filter on Newsfeed: P30Age-based Deboost on Newsfeed: P10Age-based Demotion on Newsfeed: P30Deboost on Newsfeed: P50Demotion on Newsfeed: P50EU Demotions on Newsfeed: P80Demotion on IFR: 0.50Non-rec Filter on IFR: 0.60Age gating on Newsfeed: P10ARC Demotion on Newsfeed: P50 Borderline: <ul style="list-style-type: none">Age-based Non-rec Filter on Newsfeed: P10Demotion on IFR: 0.50Non-rec Filter on IFR: P60ARC Demotion on Newsfeed: P50ARC Deboost on Newsfeed: P20
24	IG Comment	5295222083900364	06/06/2022	Bullying and Harassment = 0.95	Bullying and Harassment: <ul style="list-style-type: none">Comment filter (auto-hide): P80Comment cover: P70Advanced comment filtering: P60

Highly Confidential (Competitor)

1	2	3	4	5	6	7	8
				Borderline Hostile Speech = n/a	<ul style="list-style-type: none"> Comment downranking (demotion): P50 Preview comment filter: P50 <p>Borderline Hostile Speech:</p> <ul style="list-style-type: none"> Comment filter (auto-hide): P95 Comment cover: P80 Advanced comment filtering: P70 Comment downranking (demotion): P70 Preview comment filter: P70 		
9	IG Post	24098062513172542	11/05/2023	Bullying and Harassment = 0.95 Borderline Hostile Speech = n/a	<p>Bullying and Harassment:</p> <ul style="list-style-type: none"> Youth non-rec filter: P50 Youth demotions on connected stories: P60 Youth demotions on connected feed: P60 EU youth demotions on connected feed: P80 <p>Borderline Hostile Speech:</p> <ul style="list-style-type: none"> Youth non-rec filter: P40 		
10	11	12	13	14	15	16	17
18	Reels	7408044142611719	04/27/2024	Bullying and Harassment = 0.95 Borderline Hostile Speech = n/a	<p>Facebook</p> <p>Bullying and Harassment:</p> <ul style="list-style-type: none"> Age-based Non-rec Filter on Reels: P15 Non-rec Filter on Reels: P25 Non-rec Filter on Watch: P25 <p>Borderline Hostile Speech:</p> <ul style="list-style-type: none"> Age-based Non-rec Filter on Reels: P15 Non-rec Filter on Reels: P25 Non-rec Filter on Watch: P25 <p>Instagram</p> <p>Bullying and Harassment:</p> <ul style="list-style-type: none"> Youth demotions on connected feed: P60 		
20	21	22	23	24	25	26	27
28							

Highly Confidential (Competitor)

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28
Graphic and Violent Content																											
Graphic Violence MAD photo classifier on FB	f647163274	10/01/2024	N/A	• Mark as disturbing ("MAD"): 0.972 • Various enforcements > P50 ◦ Reels: 0.274 ◦ CFR Demotion: 0.274 ◦ CFR Deboost: 0.274 ◦ Stories demotion: 0.2529 ◦ Search filtering : 0.2529																							
Graphic Violence MAD video classifier on FB	f645517662	9/25/2024	N/A	• MAD: 0.98 • Various enforcements > P50 ◦ Reels 0.2529 ◦ CFR Demotion: 0.2529 ◦ CFR Deboost: 0.2529 ◦ Stories demotion: 0.2529 ◦ Search filtering : 0.2529																							
Graphic Violence MAD photo classifier on IG	f647151014	09/30/2024	N/A	• MAD: 0.967 • EU age gating on Feed: 0.8 • Non-EU age gating on Feed: 0.6																							
Graphic Violence MAD video classifier on IG	f639562790	09/11/2024	N/A	• MAD: 0.91 • EU age gating on Feed: 0.8 • Non-EU age gating on Feed: 0.6																							
GV Video Deletion classifier on IG	3614043688651385	02/21/2021	0.62	N/A																							
Borderline GV classifier			N/A	Facebook Youth: Various enforcements > P10 • Feed age gating: 0.03 (P10) • Feed demotion: 0.12 (P30) • Rees age gating: 0.06 (P20) • Reels age filtering: 0.09 (P25)																							

Highly Confidential (Competitor)

1	2	3	4	5	• Watch/IFR filtering: 0.09 (P25) Instagram Youth: various enforcements > P40 • Filtering: 0.341 (P50) • Demotion: 0.204 (P40) • Age gating: 0.653 (P80)
Adult Nudity & Sexual Activity					
6	Violating ANSA photo classifier on FB and IG	f644294876	09/18/2024	• 0.952 for FB • 0.96 for IG	Instagram • Age gating: 0.8 Facebook Youth filter on stories: 0.5 Youth filter on reels: 0.275 Youth filter on watch: 0.5 Youth filter on feed: 0.1
7	Violating ANSA video classifier on FB and IG	f643324421	09/16/2024	• 0.981 for FB 0.993 for IG	Instagram • Age gating: 0.8 Facebook • Youth filter on stories: 0.5 • Youth filter on reels: 0.275 • Youth filter on watch: 0.5 Youth filter on feed: 0.1
8	Violating ANSA linkshare classifier on FB and IG	f635467640	08/26/2024	0.55 for FB	Instagram • Age gating: 0.8 Facebook • Youth filter on stories: 0.5 • Youth filter on reels: 0.275 • Youth filter on watch: 0.5 • Youth filter on feed: 0.1
9	Borderline ANSA Photo Classifier on FB	f645114003	9/21/2024	N/A	• Youth filter on Reels: 0.055 • Youth filter on Watch: 0.095 • Youth filter on Feed: 0.15 • Youth filter on IFR: 0.1 • Youth filter on Story: 0.5
10	Borderline ANSA Video Classifier on FB	f645120189	9/21/2024	N/A	• Youth filter on Reels: 0.055 • Youth filter on Watch: 0.095 • Youth filter on Feed: 0.15 • Youth filter on IFR: 0.1 • Youth filter on Story: 0.5

Highly Confidential (Competitor)

1	IG Non-rec ANSA Photo Classifier on IG	f639524943	09/04/2024	N/A	<ul style="list-style-type: none"> • Age gating: 0.623 • Age filtering: 0.159 • Age demotion: 0.104
2	IG Non-rec ANSA video classifier on IG	f641128648	09/10/2024	N/A	<ul style="list-style-type: none"> • Age gating: 0.623 • Age filtering: 0.237 • Age demotion: 0.158
Child Safety					
6	CSAM classifier	FB: f588692710; IG: f565741344	08/29/2024	N/A	Facebook <ul style="list-style-type: none"> • Age filtering: P20 • Filtering for GenPop: P25 Instagram Filtering for Teens: P25
7	CSAM Solicitation classifier	FB: f636246456; IG: f566561739	09/16/2024	N/A	Facebook <ul style="list-style-type: none"> • Age filtering: P20 • Filtering for GenPop: P25 Instagram Filtering for Teens: P25
12	CSAM-S comment classifier	f636721225	08/29/2024	N/A	Facebook <ul style="list-style-type: none"> • Age filtering: P20 • Filtering for GenPop: P25 Instagram IG: no soft actions
15	Child Sexualization (CSx) classifier	f590979030	08/28/2024	N/A	Facebook <ul style="list-style-type: none"> • Age filtering: P20 • Filtering for GenPop: P25 Instagram Filtering for Teens: P25
18	CSE-I Classifier	f573679609	06/18/2024	N/A	Facebook <ul style="list-style-type: none"> • Age filtering: P20 • Filtering for GenPop: P25 Instagram IG: no soft actions
21	Child subject to objectionable content classifier	f589982602	08/20/2024	N/A	Facebook <ul style="list-style-type: none"> • Age filtering: P20 • Filtering for GenPop: P30 Instagram Filtering for Teens: P50
24	Google Content Safety API classifier	N/A	06/02/20	N/A	No soft actions